# Approach to atomistic modelling
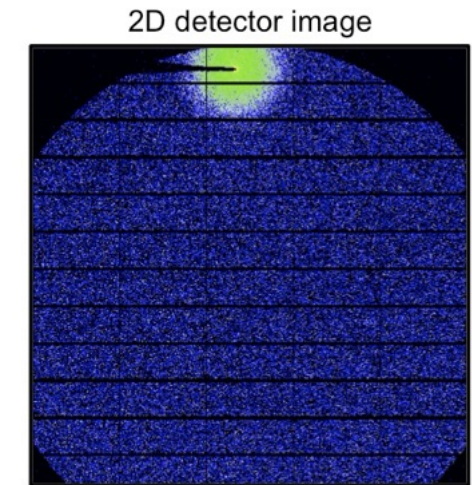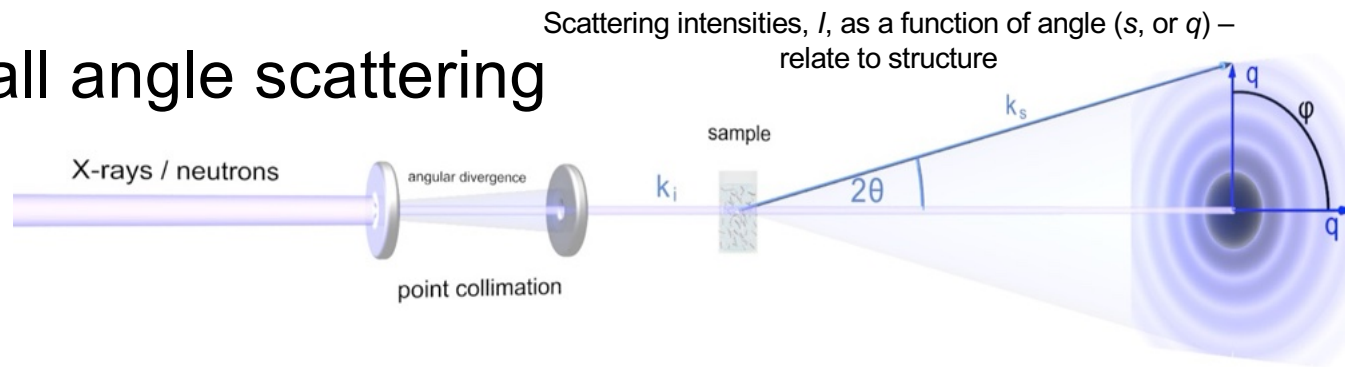# Data, quality, fitting models and modelling

Cy Jeffries

EMBL Hamburg

EMBL

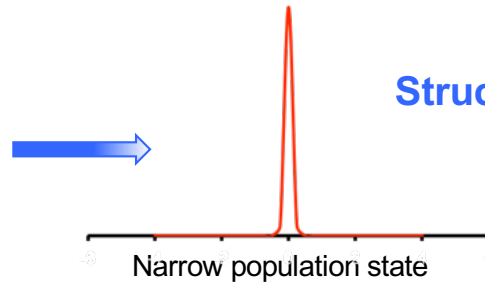# Small angle scattering

Scattering intensities, $I$, as a function of angle ($s$, or $q$) – relate to structure

X-rays / neutrons

angular divergence

point collimation

$k_i$

sample

$2\theta$

$k_s$

$q$

$\varphi$

$q$

detector

2D detector image

**Sample Properties:** <u>macromolecules and conjugates in solution</u>

*always illuminate a population*

Sample

Monodisperse

$$I(s)=\sum_k v_k I_k(s)$$

Narrow population state

**Structurally homogeneous**

**Single-state model**

Polydisperse
Different conformations
Different oligomeric states

population state(s) distribution

**Structurally heterogeneous**

Multi-state models

**multiple population states:mixtures**

*Temperature*
*Time*
*Chemical environment*

population state transition

EMBL

# Small angle scattering

## Wave properties: Scattering amplitudes

The magnitude of the **coherent scattering amplitudes** as a function of angle relates to spatial correlations between scattering centres.
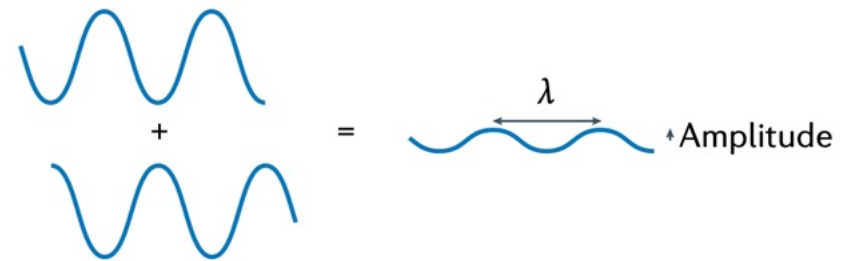


Incident beam
$\vec{k_i}$
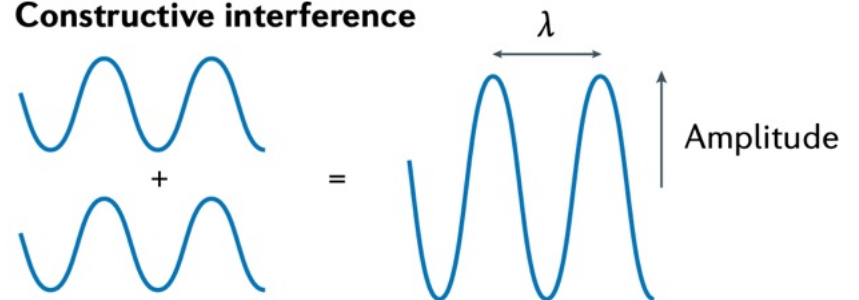$\lambda$
Amplitude
Scattering wave amplitudes

When the incident beam is **scattered elastically** (no change in $\lambda$) and if preserved distance correlations exist between the scattering centres, a **coherent wavefront** develops that emanates from the sample where both constructive and destructive interference occurs in the wave amplitudes.

**Destructive interference**



$\lambda$

+ = Amplitude

**Constructive interference**



$\lambda$

+ = Amplitude

EMBL

Elastic coherent small-angle scattering =
'Sum the waves game.'



'Waves cancel' (amplitudes cancel)

Line of destructive interference

Amplitude pattern across the coherent wave front relates to the correlated distance between the two scattering centres

**Preserved distance** correlation between two scattering centres.

Line of constructive interference

'Waves add' (amplitudes add)

Coherent wave front
elastic scattering – no wavelength change

Of course, macromolecules have many, many atom pair distance correlations within extent of their volume boundary. The coherent wave front is derived from the **sum** of the scattered waves from all of these correlations as a function of angle.

EMBL

# More formally:

If the distances, *r*, between the atoms of a marcomolecule are preserved then the amplitudes of the *coherent* wave front through *s* are proportionate to the sum of the atomic scattering factors (i.e., probability to scatter) weighted by the *distribution of the distances* between scattering pairs.

$$s = \frac{4\pi \sin\theta}{\lambda}$$

!! *s* can be defined in a number of ways!!

$Q = q = h = \mu = k = s$

Sometimes; $S = 2\sin\theta/\lambda = 2\pi s$

$$A(s) = \sum_{i=1}^{N} b_i \left( e^{i\vec{s}\cdot\vec{r}} \right)$$
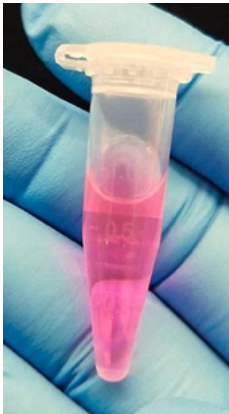
Spherical wave bit

'Scattering factor': relates to the atomic cross section, i.e., scattering length, or probability of an atom to scatter for every atom in the sample.

EMBL

# The issue?

We cannot access the amplitudes experimentally. We measure the *intensity* of the scattered radiation, $I(s)$.

For solution-based SAXS, the sample particles are tumbling in solution!
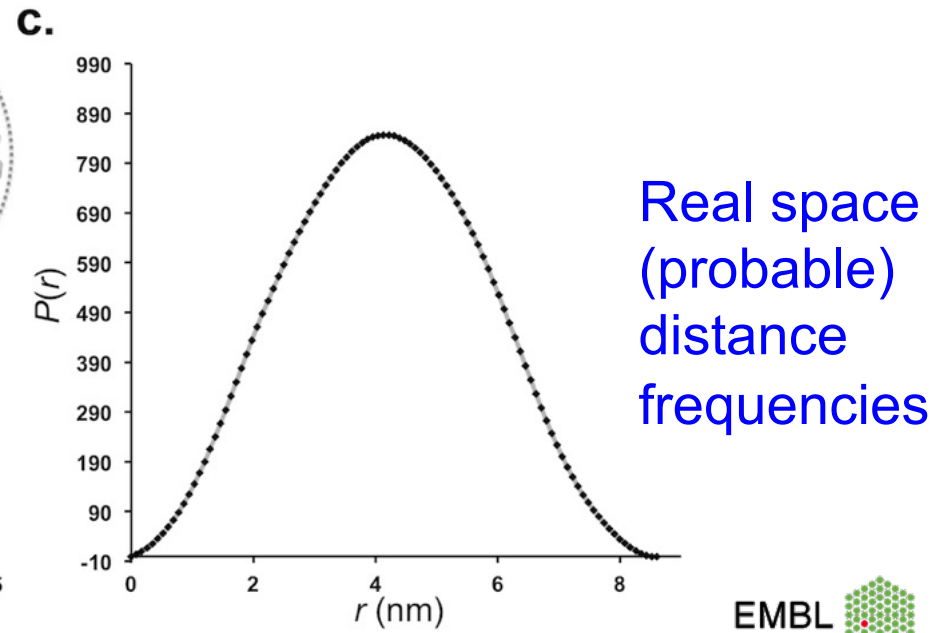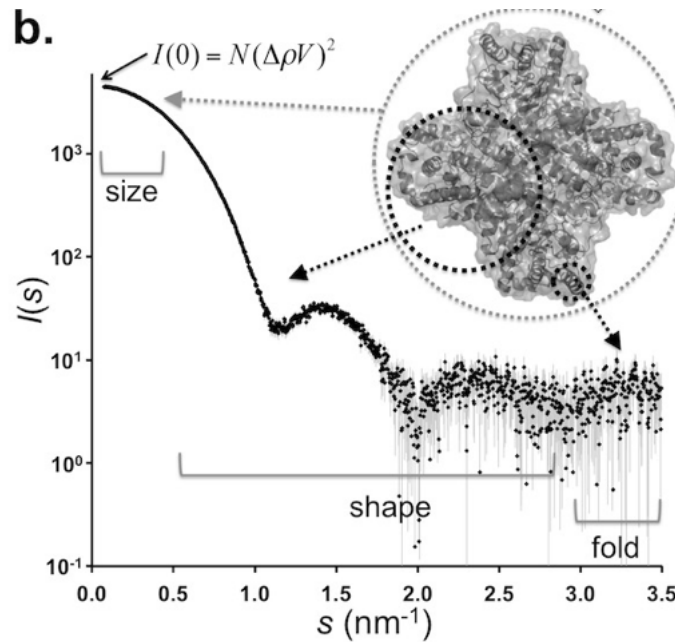
Sample



$$I(s) = \left\langle A(s)A(s)^* \right\rangle$$

⟶ All orientations considered, i.e., isotropic scattering

$I(s)$ fundamentally boils down to the form factor of the particles, $P(s)$, their volume and scattering power!

EMBL

The scattering intensity $I(s)$ – and thus the associated form factor in reciprocal space – relates to an atom-pair distance distribution function of the particle $p(r)$ in real space by a Fourier transform:

$$I(s) = 4\pi \int_0^{D_{max}} p(r) \frac{\sin(sr)}{sr} dr \qquad p(r) = \frac{r^2}{2\pi} \int_0^\infty s^2 I(s) \frac{\sin(sr)}{sr} ds$$
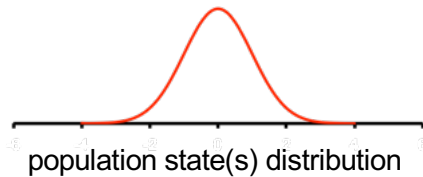


Reciprocal space intensity

Real space (probable) distance frequencies

EMBL

# How does the sample properties combined with the measurement impact our approach to modelling data?
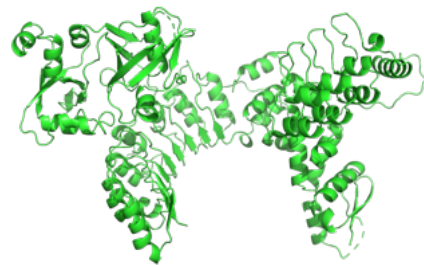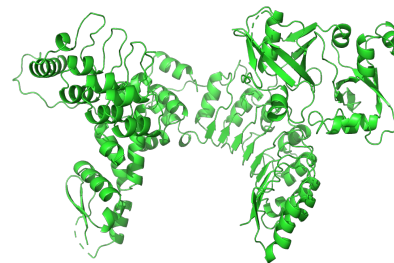
Populations of macromolecules in the sample

Sum of spherically averaged scattering amplitudes

population state(s) distribution

$$I(s) = \left\langle A(s)A(s)^* \right\rangle$$

Distance distributions are encoded in the scattering intensities – not x,y,z atomic coordinates

L-amino acids

D-amino acids

EMBL

# The Complication

For biological macromolecules in solution, *we forgot the solution!*

It is obvious that macromolecules of a sample will scatter. The amplitudes arising from preserved distance correlations will **sum** to produce coherent scattering intensities at low angle.

The lower the angle (lower $s$), the *longer* the correlated distances, $d$:

$$s = 2\pi/d$$

*However*, the solution, i.e., the solvent of the sample, also scatters! As the solvent (hopefully) does not have any time-preserved long-range distance correlations, its scattering contributions add as a 'flat incoherent background' in the SAS regime.
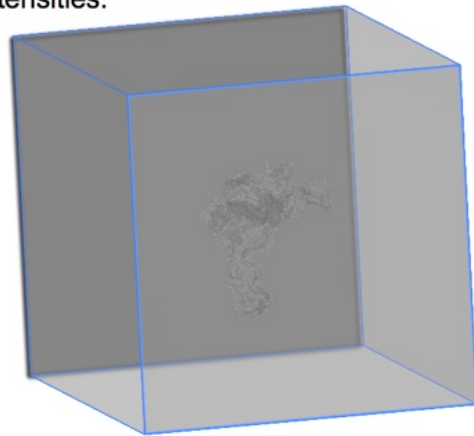
EMBL

# But you need contrast as well: 'excess scattering power'

$I(s)$ in the small-angle region depends, and indeed only arises, if there is a difference between the average scattering length density of the **solvent** and the average scattering length density of the <u>particles of interest</u>. This difference is known as ***contrast*** and is represented as
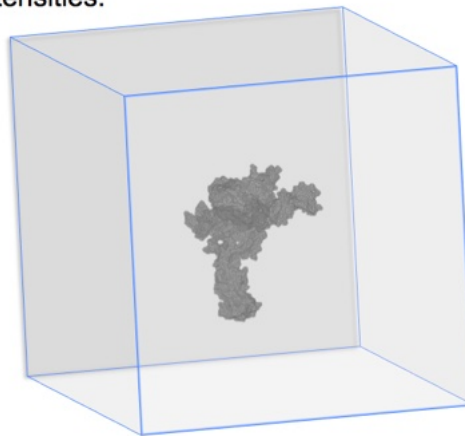
$$\Delta\rho = \overline{\rho} - \rho_{s},$$

where $\overline{\rho}$ and $\rho_{s}$ are the mean scattering length densities of the particle and the solvent, respectively.

Low contrast = weaker coherent scattering intensities.

High contrast = stronger coherent scattering intensities.

$$I(s) \propto \Delta\rho^{2}$$

EMBL

# How do I calculate the Contrast

http://smb-research.smb.usyd.edu.au/NCVWeb/

**MULCh**: Modules for the analysis of small-angle neutron contrast variation data from bio-molecular assemblies.

**For X-rays**: Convert the SLD, $\rho$ ($10^{10}$ cm$^{-2}$) to electron density by dividing by the Thomson electron radius: $2.8179 \times 10^{-13}$ cm. The answer is in e/cm$^3$, so divide again by $10^{24}$ to get e/Å$^3$…or more quickly:

$$\frac{\rho}{28.179} \quad \text{e/Å}^3$$

11          01/10/2024



A: Define the solvent

ModULes For The Analysis Of Small-angle Neutron Contrast Variation Data From Bio-molecular Assemblies

Contrast: Module For Estimating The Contrast Of Bio-molecular Assemblies

Upload an existing input file:
Upload Contrast File     (Upload txt input if avaialble)

Project Title: VH Ab lysozyme          i) Input title of project

Number of contrast points: 0

Number dissolved species in the solvent: 3          ii) # molecules in solvent = 3

iii) For small molecules: input atomic formula and concentrations

P = protein; D = DNA; R = RNA; M = molecule

B: Define macromolecules          i) # components in subunit 1 = 1

ii) Choose level of deuteration.

iii) Amino acid sequence

iv) # components in subunit 2 = 2
v) Choose level of deuteration

vi) Amino acid sequence

vii) Bound calcium; 2 per subunit.

C: $\rho$ and $\Delta\rho$ output

| | $\rho$ ($10^{10}$cm$^{-2}$) | | | $\Delta\rho$ ($10^{10}$cm$^{-2}$) | | |
| | 1 | 2 | Solvent | 1 | 2 | Total |
|---|---|---|---|---|---|---|
| X-RAY | 12.515 | 12.580 | 9.454 | 3.061 | 3.127 | 3.095 |
| NEUTRON Fraction $^2$H$_2$O in solvent, (SANS). 0.0 | 1.957 | 4.579 | -0.545 | 2.502 | 5.123 | 3.869 |
| 0.1 | 2.112 | 4.712 | 0.146 | 1.967 | 4.566 | 3.322 |
| 0.2 | 2.268 | 4.844 | 0.836 | 1.432 | 4.008 | 2.775 |
| 0.3 | 2.423 | 4.977 | 1.526 | 0.897 | 3.451 | 2.228 |
| 0.4 | 2.579 | 5.110 | 2.217 | 0.362 | 2.893 | 1.682 |
| 0.5 | 2.734 | 5.242 | 2.907 | -0.173 | 2.335 | 1.135 |
| 0.6 | 2.890 | 5.375 | 3.597 | -0.707 | 1.778 | 0.588 |
| 0.7 | 3.045 | 5.508 | 4.288 | -1.242 | 1.220 | 0.042 |
| 0.8 | 3.201 | 5.641 | 4.978 | -1.777 | 0.663 | -0.505 |
| 0.9 | 3.356 | 5.773 | 5.668 | -2.312 | 0.105 | -1.052 |
| 1.0 | 3.512 | 5.906 | 6.359 | -2.847 | -0.453 | -1.598 |
| Calculated match-point ($f_{D_2O}$) | | | | 0.468 | 0.919 | 0.708 |

i) X-ray contrast for SAXS

ii) Total neutron contrast for SANS

Individual component SANS matchpoints: v/v $^2$H$_2$O.

Whole complex SANS matchpoint: v/v $^2$H$_2$O.

## *After* background subtraction…

*I(s)* will represent the time and rotationally averaged squared scattering amplitudes from the particle population expressed as the summed contribution from each individual particle, *i*, in the sample.

The scattering intensity

Is the SUM of all macromolecules averaged over all orientations.

The structure factor or 'between particle' contributions

$$I(s) = \sum_{i}^{n} [(\Delta\rho_i V_i)^2 P_i(s)] S(s)$$

Weighted by the contrast and volume SQUARED of all macromolecules

The form factor of all macromolecules within the sample

EMBL

For a PURE, MONODISPERSE and IDEAL sample

The **concentration**.

$$I(s) = N(\Delta\rho V)^2 P(s)$$

If all particles are identical, and do not interact, the *I(s)* profile (*after* background scattering has been subtracted) will represent the time and rotationally averaged squared scattering amplitudes, i.e., the scattering intensity, from a **SINGLE PARTICLE**.

EMBL

# How do I *maybe* know I have an ideal system?
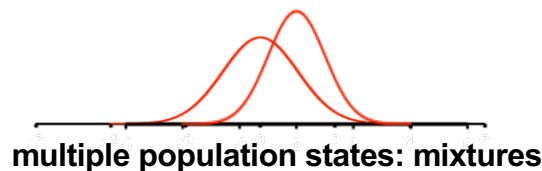
The molecular mass estimates through a concentration series.

*The MM, the MM, the MM, the MM, the MM, the MM.*

(+/-10 %)

***Think about this*** – there is no point generating a single model to describe a 100 kDa protein if the experimental MW of the protein from SAS is 125 kDa (probably a mixture).
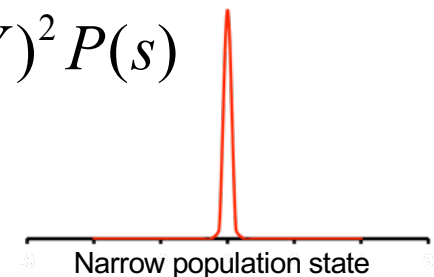
$$I(s) = \sum_i^n \left[ (\Delta\rho_i V_i)^2 P_i(s) \right]$$

$$I(s) = N(\Delta\rho V)^2 P(s)$$

If *i* are not identical, model as a mixture

**multiple population states: mixtures**

If *i* are all identical, model as a single particle

Narrow population state

EMBL

The case of intrinsically disordered proteins or modular proteins connected by flexible linkers.

$$I(s) = \sum_{i}^{n} [(\Delta\rho_i V_i)^2 P_i(s)]$$

Model as a structural ensemble!



population state distribution

**MM is correct!**

*BUT*

Still not ideal, i.e., cannot be modelled as a non-interacting single particle because the protein is structurally heterogeneous!

EMBL

# *I*(0)

At zero angle (*s* = 0) the magnitude of *I*(*s*) will primarily depend on the number of scattering centres within the bound squared-volume of a macromolecule – independent of the shape – weighted by the concentration and contrast squared:

$$I(0) \approx N(\Delta \rho V)^2$$

From this parameter, it is possible to obtain the ***molecular weight***.

Data scaled to a standard protein with a KNOWN concentration and molecular weight

$$\mathrm{MW}_{sample} = \frac{I(0)_{sample}\, \mathrm{N_A}}{c_{sample}(\Delta \rho \upsilon_{sample})^2}$$

$$\mathrm{MW}_{sample} = \frac{I(0)_{sample}}{I(0)_{standard}} \times \frac{c_{standard}\,\Delta\rho^2_{standard}}{c_{sample}\,\Delta\rho^2_{sample}} \times MW_{standard}$$

Absolute scaling - requires partial specific volume and contrast.

An assumption that a target has a similar scattering length density and partial specific volume as the secondary standard!

EMBL

# Porod volumes and Kratky plot

The determination of MW from $I(0)$ requires an accurate assessment of the concentration of a protein in solution that in and of itself can be difficult to determine!

An alternative concentration-independent estimate of MW is based on the volume of a protein in solution. Porod showed that for uniform particles with sharp boundaries the excluded volume $Vp$ can be calculated as:

$$V_p = \frac{2\pi^2 I(0)}{Q}$$

where $Q$ is the Porod invariant or the area under a plot of $I(s)s^2$ vs $s$ calculated to s = ∞, or Kratky plot.

The $Vp$ of a protein in $nm^3$ is typically 1.5–1.6 times the MW in kilodaltons (kDa).

However, caution must be applied when dealing with highly anisotropic or highly flexible/disordered proteins. In the case of flexible, or rod-like proteins, the decay in scattering intensities at high angle deviates sufficiently from Porod's law that the estimation of $Q$ will incur errors in the volume estimation!

EMBL

# Useful ATSAS tools

ATSAS tool: *datporod*     ATSAS tool: *datmow*     ATSAS tool: *datvc*

At the command prompt (.cmd, terminal, etc) type:

datporod  filename.out     datmow  filename.out     datvc  filename.out
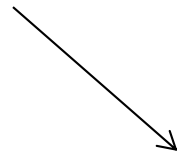
↓     ↓     ↓

Porod volume estimate.
For proteins, convert to MM by
dividing by 1.5-1.6
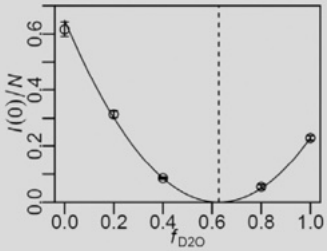
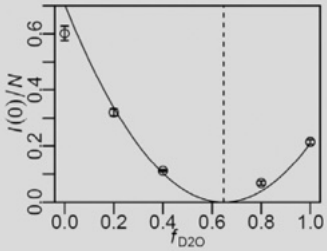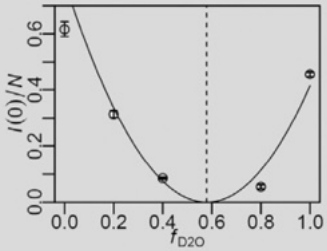MM estimate of proteins using
the method of Fischer et al.
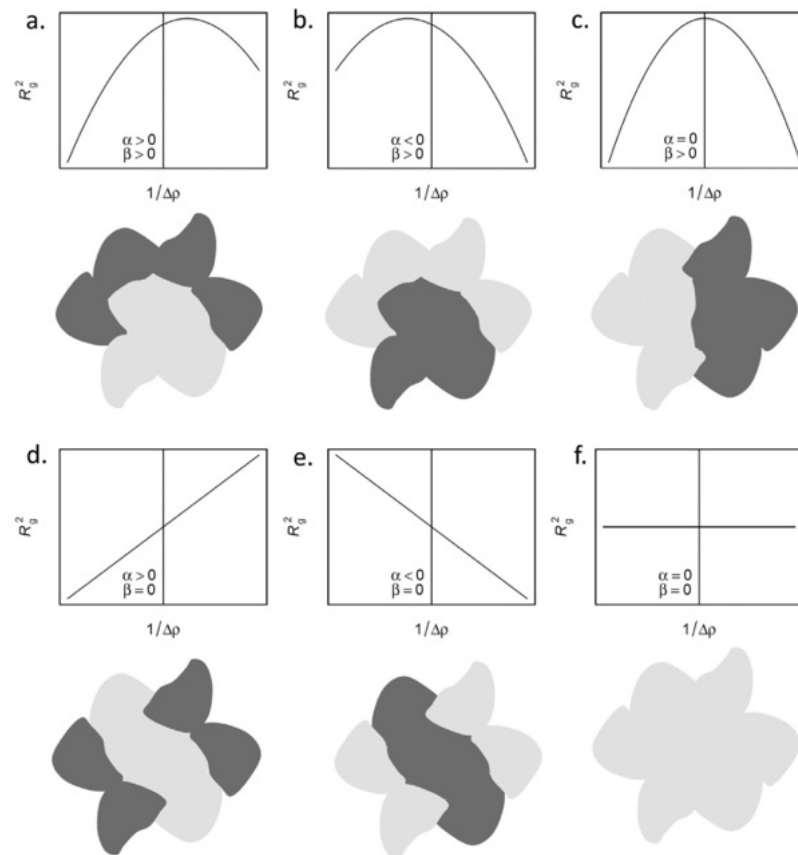*SAXMOW*

MM estimate of proteins using the
method of Rambo and Tainer.

ATSAS tool: *datmw*

EMBL

# Why is all this stuff important for SANS?

$$R^2_{g,obs} = R^2_C + \frac{\alpha}{\Delta\bar{\rho}} - \frac{\beta}{\Delta\bar{\rho}^2},$$

# Modelling SAS data – before you leap into danger.

- Understand the data – get the unit right, nm or Å, etc.

- Extract structural parameters and additional information *BEFORE* you begin modelling: if there is one thing you can trust it is the structural parameters from SAS data!

| | |
|---|---|
| **Part 1 of your validation toolbox** $\longrightarrow$ | <ul><li>Radius of gyration ($R_g$) maximum particle dimension ($D_{max}$), volume ($V$).</li><li>Molecular mass estimates ($MM$).</li><li>Probable frequency of distances ($r$) within single particles ($p(r)$ vs $r$), i.e., *global* shape and structural information.</li><li>*Scaling parameters* – compact, flexible, flat, rod, hollow.</li><li>*Useful data range!*</li><li>*The AMBIGIUTY of the data!*</li><li>Size distributions and volume fractions.</li></ul> |

EMBL

# Modelling SAS data – before you leap into danger.

- Obtain as much information as possible about your system.
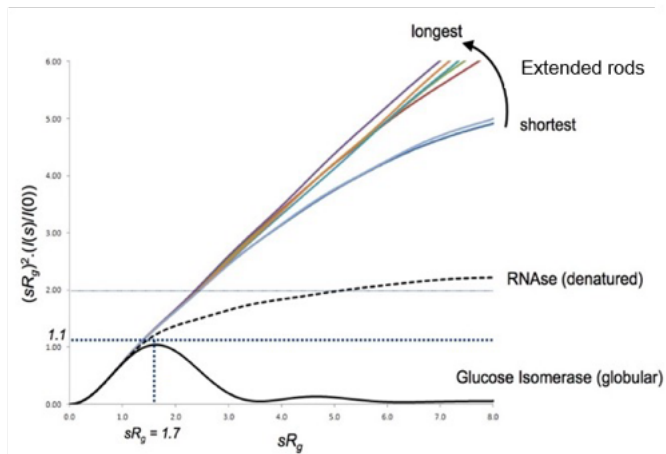
**Part 2 of your validation toolbox** →

- For example, obtain the *EXACT* sequence of the protein(s), RNA, DNA, glycans, etc actually used for the SAS measurement. *ALL atoms scatter*, so you have to take into account *ALL of the mass* in your modelling!

- Obtain the *CORRECT PDB (or cif)* files (i.e., atomic coordinate files). *ALL atoms scatter*, so you have to take into account *ALL of the mass* in your modelling!

- If required, calculate the *CONTRAST* of your system – especially important for neutrons; (on occasion, for SAXS, convert to electron density difference.)

- Obtain restraints derived from complementary methods – in particular *CONTACT* information (e.g., from NMR, cross-linking mass-spectrometry, FRET, bioinformatics, Alphafold.)

- Know the *STOICHIOMETRY* and from this, the estimated *SYMMETRY*. Obtain the MM estimate from SAS or other methods, e.g., MALLS, AUC, mass-photometry, etc.

EMBL

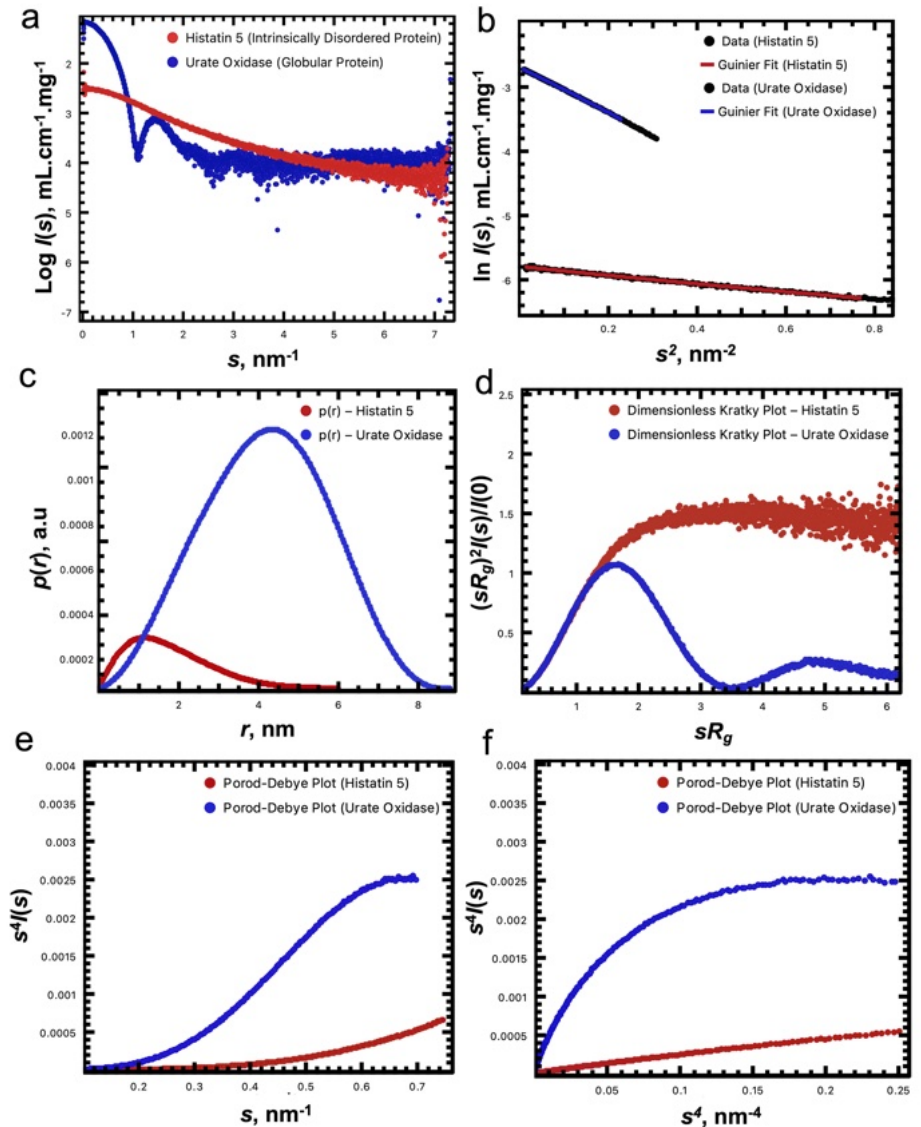# Get to *know your data* before you model!

Fundamental plots

1) Guinier
2) *P(r)*
3) Dimensionless Kratky
4) Porod-Debye



Receveur-Bréchot & Durand (2012) *Current Protein and Peptide Science*, 13, 55-75.
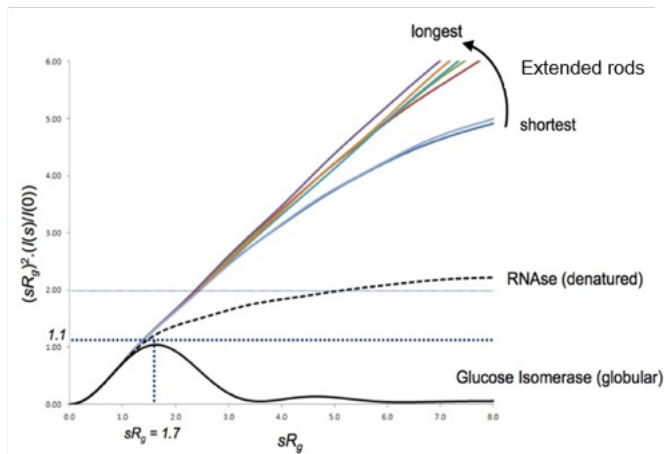Durand D, Vivès C, Cannella D, Pérez J, Pebay-Peyroula E, Vachette P, Fieschi F. (2010) *J Struct Biol.* 169: 45-53.

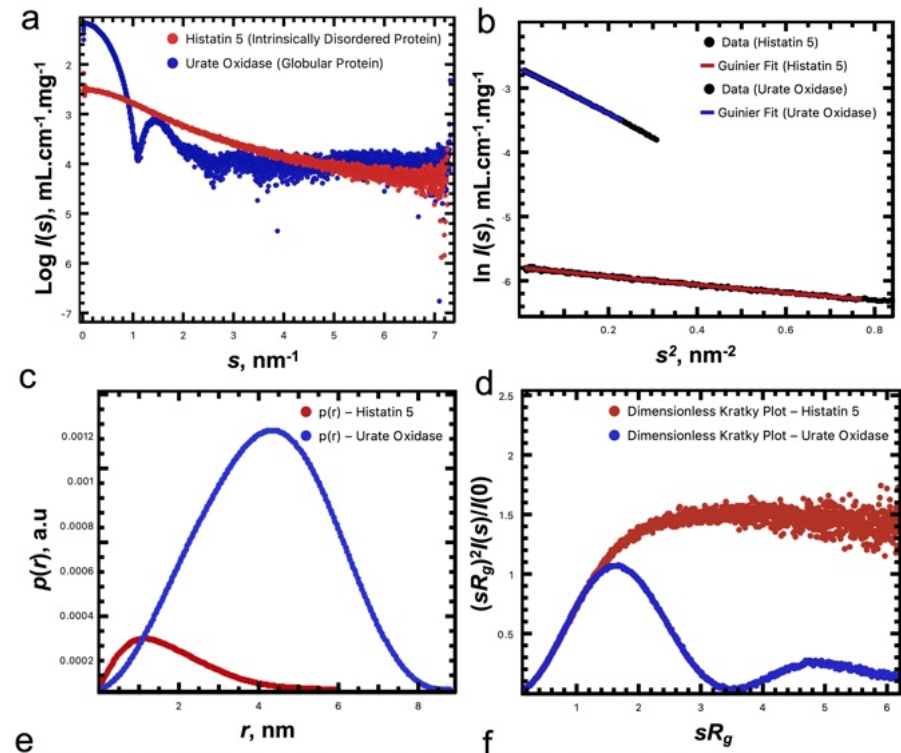# Get to *know your data* before you model!

Fundamental plots

1) Guinier
2) P(r)
3) Dimensionless Kratky
4) Porod-Debye



Receveur-Bréchot & Durand (2012) *Current Protein and Peptide Science*, 13, 55-75.
Durand D, Vivès C, Cannella D, Pérez J, Pebay-Peyroula E, Vachette P, Fieschi F. (2010) *J Struct Biol.* 169: 45-53.

# Get to *know your data* before you model!

2) Five programs:

- AutoRg – for first assessments of $s_{min}$

- SHANUM – define the useful $s_{max}$.

- DATCLASS – machine-learning method for the rapid geometric classification of SAXS data (from proteins).

- DARA – kd-tree searching of the PDB + Alphafold DB for similar scattering profiles.

- AMBIMETER – assess the ambiguity of the scattering data.

EMBL
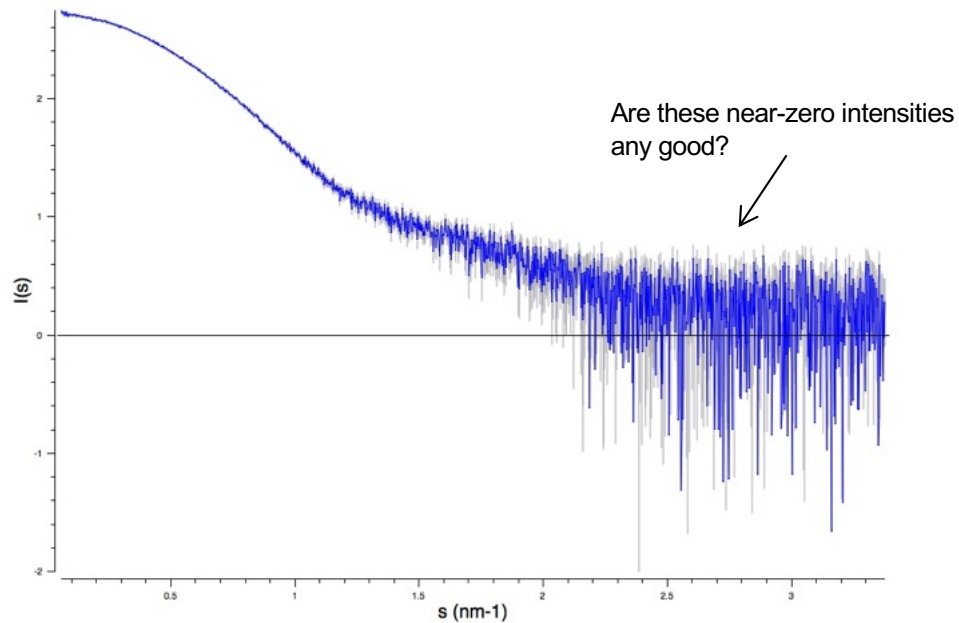
# AutoRg and $s_{min}$



Near beam-stop garbage

After cleaning up! **Is there sufficient data at the very lowest $s$ to encompass the particle size?**
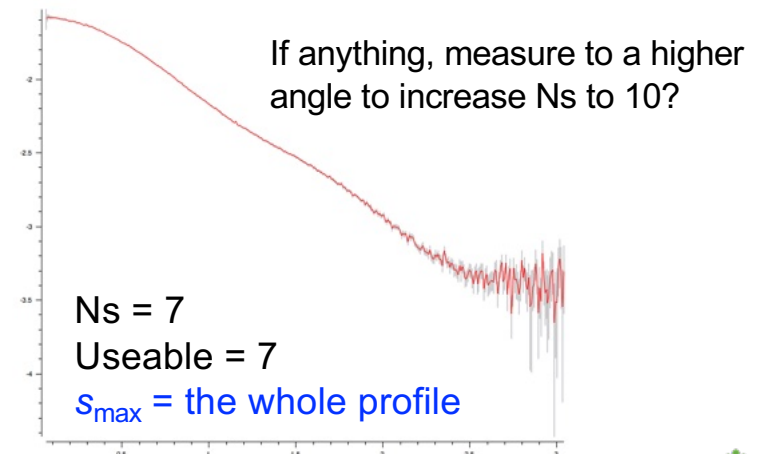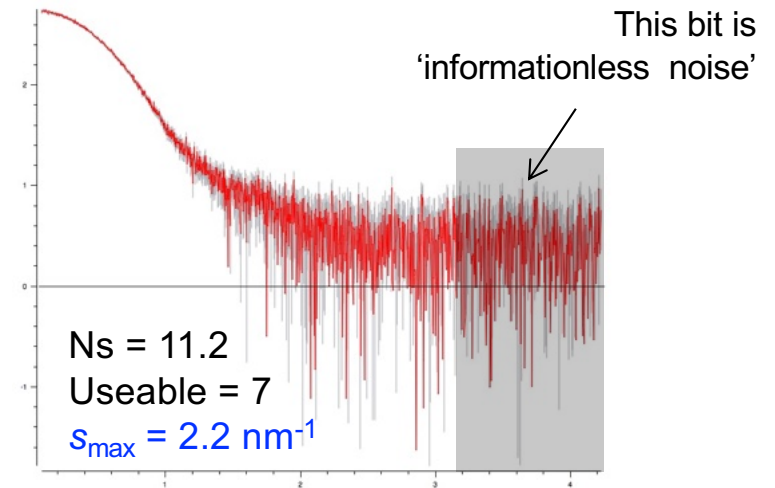
At minimum, $s_{min}$ should $= \pi/D_{max}$
Better rule of thumb, $s_{min} = 1/D_{max}$

EMBL

# Shanum and $s_{max}$

Shanum takes into account the statistical variance in the data to assess the useable $s_{max}$.

Are these near-zero intensities any good?

This bit is 'informationless noise'

Ns = 11.2
Useable = 7
$s_{max}$ = 2.2 nm$^{-1}$

If anything, measure to a higher angle to increase Ns to 10?

Ns = 7
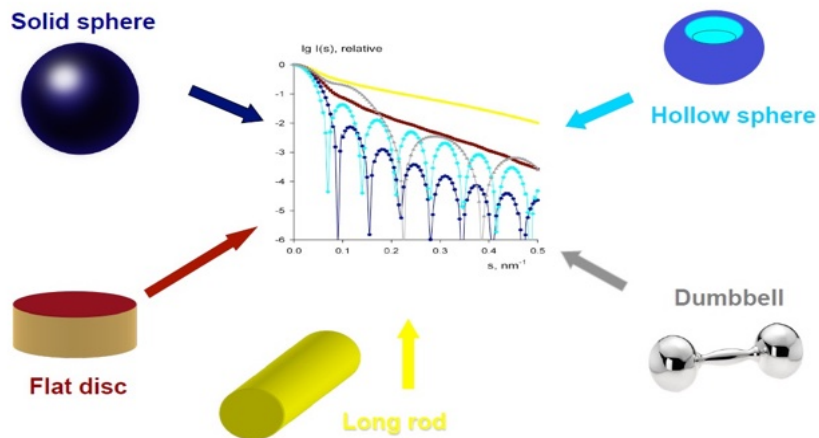Useable = 7
$s_{max}$ = the whole profile

Shanum will also estimate $D_{max}$ for you
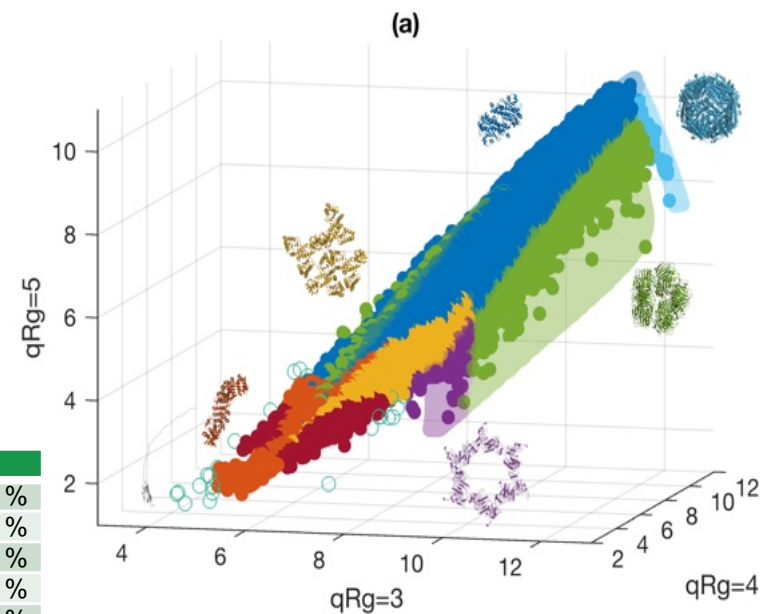(or you can enter it yourself)

EMBL

# Datclass

- Classification of a protein shape using machine learning methods based on the scattering profiles calculated from a continuum of 488 000 geometric objects **including intrinsically disordered polymers**
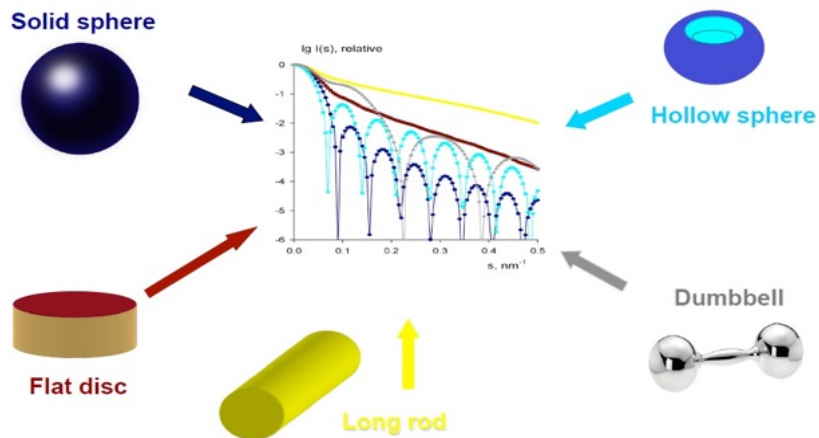
99.98% of the PDB maps into the classifier space.



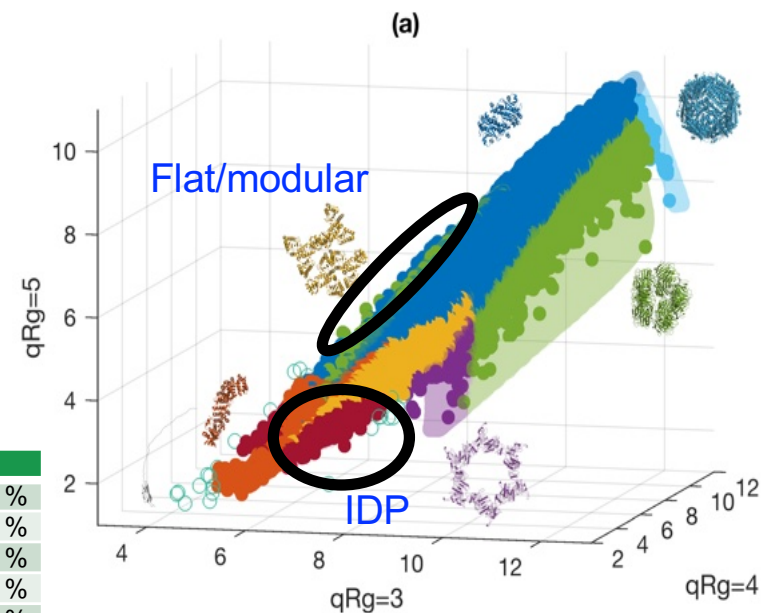| Class Label | PDB | |
|---|---|---|
| Unknown | 25 | 0. 02 % |
| Compact | 122.913 | 74.05 % |
| Extended | 5.382 | 3.24 % |
| Flat | 9.734 | 5.86 % |
| Ring | 154 | 0.09 % |
| Compact hollow | 26.909 | 16.21 % |
| Hollow sphere | 125 | 0.08 % |
| Random Chain | 740 | 0.45 % |
| Total | 165.982 | 100.00 % |

EMBL

# Datclass

- Classification of a protein shape using machine learning methods based on the scattering profiles calculated from a continuum of 488 000 geometric objects **including intrinsically disordered polymers**

99.98% of the PDB maps into the classifier space.



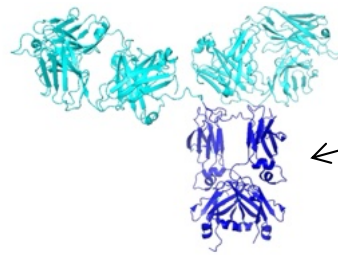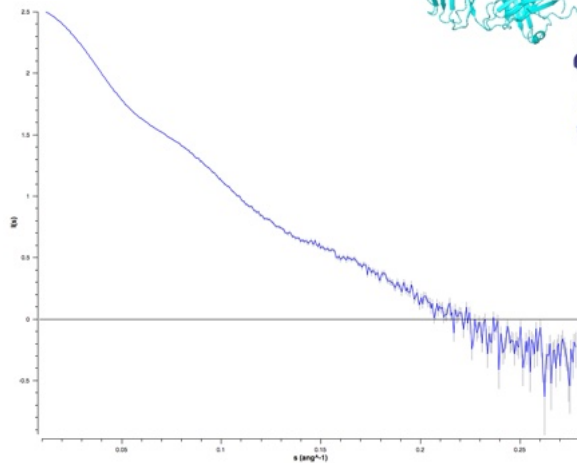| Class Label | PDB | |
|---|---|---|
| Unknown | 25 | 0. 02 % |
| Compact | 122.913 | 74.05 % |
| Extended | 5.382 | 3.24 % |
| Flat | 9.734 | 5.86 % |
| Ring | 154 | 0.09 % |
| Compact hollow | 26.909 | 16.21 % |
| Hollow sphere | 125 | 0.08 % |
| Random Chain | 740 | 0.45 % |
| Total | 165.982 | 100.00 % |

For IDP and Flat/modular = an ensemble approach might be considered!

EMBL

# DARA

Kd-tree nearest neighbour search of a .dat file or GNOM.out file against calculated SAXS profiles – PDB and Alphafold.

https://dara.embl-hamburg.de/

Combine DARA output with secondary structure prediction (predicted all β-strand). E.g., YAPSIN:
http://www.ibi.vu.nl/programs/yaspinwww/
E.g., ProteinPredict

IgG or IgA like scattering

MW estimates!



DARA neighbours

| | Fit | $x^2$ | PDB ID | Download model | MW | Volume | $R_g$ | $D_{max}$ |
|---|---|---|---|---|---|---|---|---|
| 1 | | 13.80 | ⊞ 1HZH | 6% α  45% β | 150.1 kDa | 236 nm³ | 5.3 nm | 17.3 nm |
| 2 | | 35.47 | 1R70 | 0% α  0% β | 148.7 kDa | 336 nm³ | 5.2 nm | 15.1 nm |
| 3 | | 39.65 | ⊞ 3K1M | 45% α  21% β | 139.9 kDa | 225 nm³ | 4.9 nm | 16.1 nm |
| 4 | | 49.64 | ⊞ 2FFL | 48% α  12% β | 160.5 kDa | 255 nm³ | 5.2 nm | 18.9 nm |

EMBL

# Ambimeter

Based on a set of (several thousand) shape topologies with pre calculated scattering profiles.
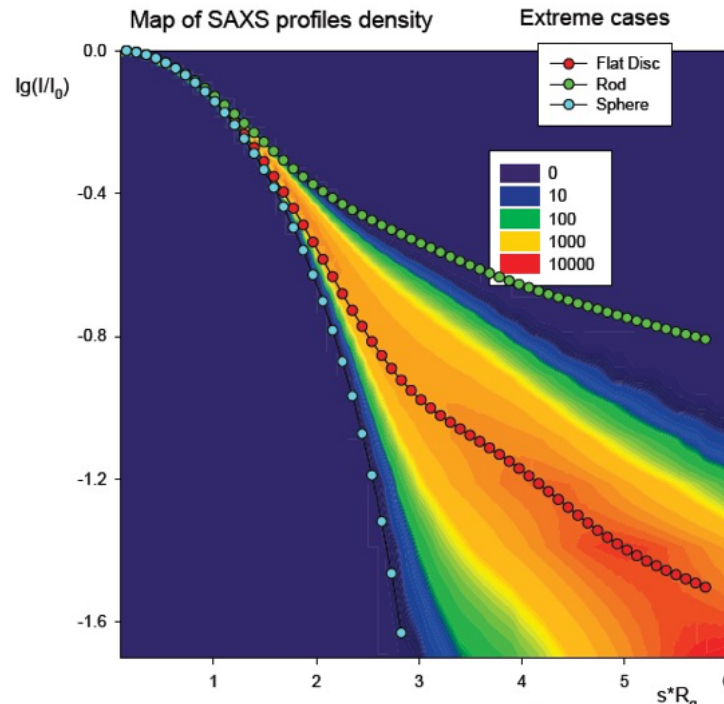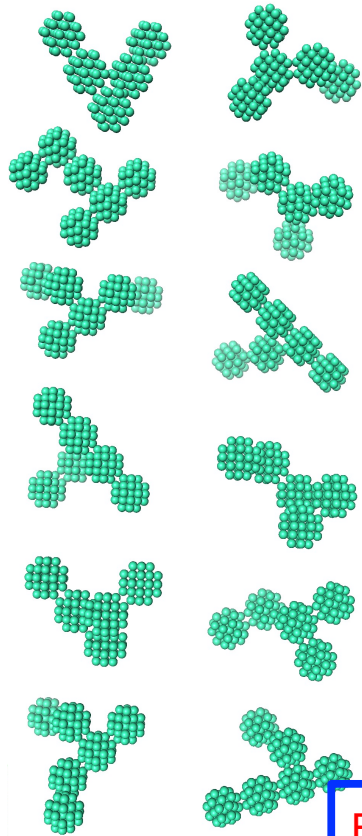


- Provides a sense of how ambiguous a dataset is with respect to fitting models.
- An ambimeter score of 0 to 1.5(ish) are 'potentially unique' shapes.
- An ambimeter score of 2.9, for example, is very highly ambiguous.

What to do?
- Always run modelling routines several times!
- Use information from other techniques.
- Perform parallel modelling against several SAS datasets.

EMBL

# Ambimeter

Based on a set of (several thousand) shape topologies with pre calculated scattering profiles.
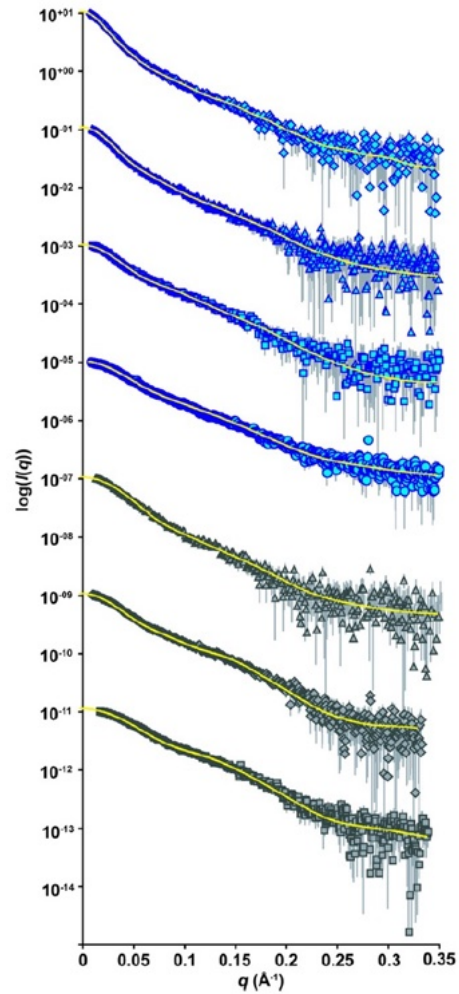


Map of SAXS profiles density — Extreme cases

- Provides a sense of how ambiguous a dataset is with respect to fitting models.
- An ambimeter score of 0 to 1.5(ish) are 'potentially unique' shapes.
- An ambimeter score of 2.9, for example, is very highly ambiguous.
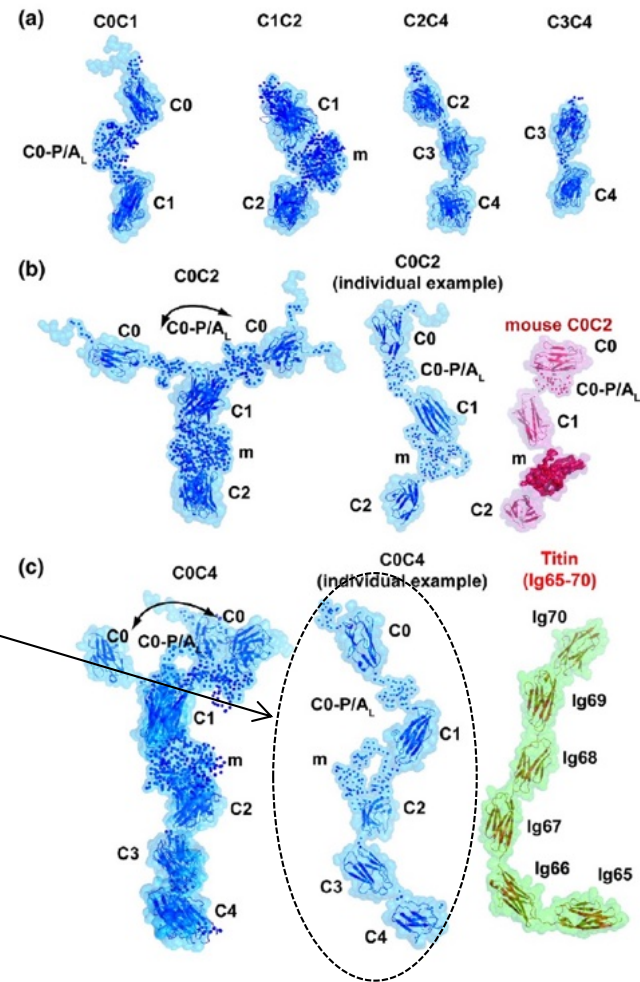
What to do?

- Always run modelling routines several times!
- Use information from other techniques.
- Perform parallel modelling against several SAS datasets.

Flat things are highly ambiguous. Classification as 'flat/modular' – are these modular with flexible linkers or just highly ambiguous?

EMBL

# More information = less ambiguity!



Parallel SAXS modelling of domain and domain constructs (truncation mutants)

Final Target: fits but importantly DOES NOT describe *in toto* what is going on

EMBL

# More information = less ambiguity! From other methods.

Xay-ray crystallography
NMR
FRET
Electron microscopy
Mass-spec (HDX, cross linking)
Predictive methods (Alphafold)

EMBL

# More information = less ambiguity! From other methods.

Xay-ray crystallography
NMR
FRET
Electron microscopy
Mass-spec (HDX, cross linking)
Predictive methods (Alphafold)



STRUCTURAL
BIOLOGY

ISSN 2059-7983

research papers

Validation of ele...scopy maps using
solution sm......y scattering

Kristi... .. Pedersen*

..., and Interdisciplinary Nanoscience Center (iNANO), Aarhus University, Gustav Wieds Vej 14,
...nark. *Correspondence e-mail: jsp@chem.au.dk

Acta Cryst. (2024). D8... https://doi.org/10.1107/S2059798324005497

*Wait for Sergei Grudinin and Dina Schneidman-Duhovny lectures Today! Atomistic modelling*
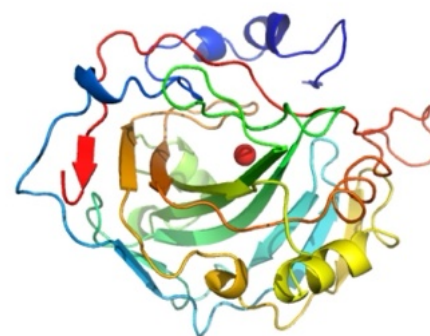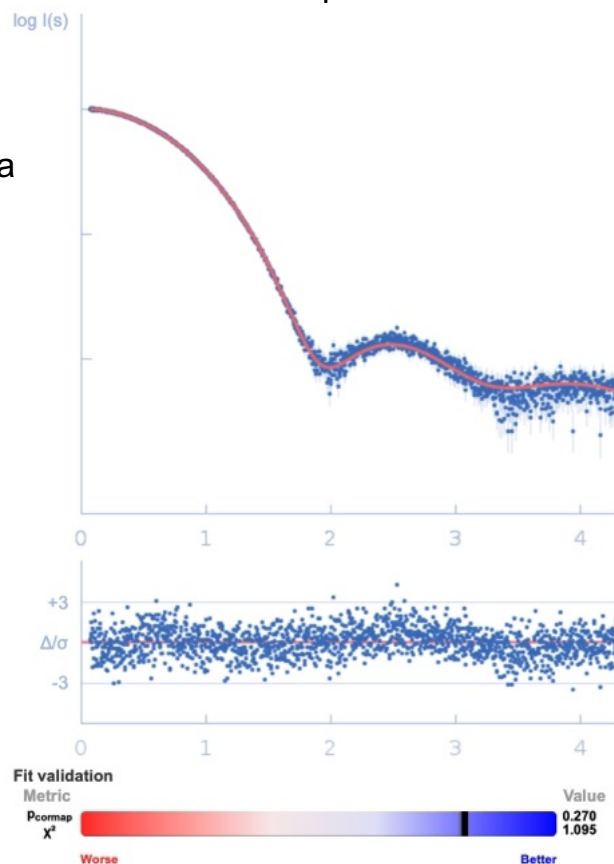
EMBL

# Summary: Know your data.

- $R_g$ and $I(0)$ from Guinier and $p(r)$ – check for consistency through a concentration series. Identify concentration independent interparticle interactions: coulombic-repulsive or aggregation. Deal with it.

- *Prepare your data for modelling:* $s_{min}$ and $s_{max}$ from AUTORG and SHANUM. Make sure $s_{min}$ is at least = $\pi/D_{max}$, or better yet, $1/D_{max}$!

- Molecular mass estimates – very important for guiding the modelling approach!

- Ambiguity.

EMBL

Any questions so far?

EMBL

# Lets do some atomistic model fitting!

https://www.sasbdb.org/data/SASDFP8/

Assessing SAS-data model fit

Atomistic all-atom model – e.g., from the PDB, Alphafold, etc

Error-normalized residual plot– looking systematic deviations between the calculated model scattering intensities and the experimental data. CORRECT ERRORS.

CorMap $P$ and $x^2$ evaluations

# Assessing SAS-data model fits – Methods abound!

| Approach | Modeling of the hydration layer | Representation of the molecule | References |
|---|---|---|---|
| CRYSOL | Implicit layer using an envelope function | All-atom | Svergun et al. *J. Appl. Cryst.* (1995) |
| AXES | Explicit water molecules using equilibrated water boxes | All-atom | Grishaev *et al. JACS* (2010) |
| FoXS | Implicit layer based on surface accessibility | All-atom or coarse-grained | Schneidman-Duhovny *et al. NAR* (2010) |
| HyPred | Explicit water molecules based on MD simulations | All-atom | Virtanen *et al. Biophys. J.* (2011) |
| AquaSAXS | Solvent-density map using the dipolar PB-Langevin approach | All-atom | Poitevin *et al. NAR* (2011) |

CRYSON – for SANS

FoXS – Debye formula

$$I(q) = \sum_i \sum_j w_i w_j \frac{\sin q r_{ij}}{q r_{ij}}$$

**Dina Schneidman-Duhovny…is here!**

EMBL

# Assessing SAS-data model fits – Methods abound!



| Approach | Modeling of the hydration layer | Representation | References |
|---|---|---|---|
| CRYSOL | Implicit layer us... envelope funct... | | |
| AXES | Explicit water m... using equilib... boxes | | |

**WAXIS – molecular dynamics for the hydration layer!**

**PEPSI-SAXS and PEPSI-SANS**

# Assessing SAS-data model fits – Methods abound!

| Approach | Modeling of the hydration lay... | Representation | References |
|---|---|---|---|
| CRYSOL | Implicit layer us... envelope funct... | | |
| AXES | Explicit water m... using equilib... boxes | | |

**Jochen Hub…is here!**

**WAXSIS – molecular dynamics for the hydration layer!**

**Sergei Grudinin…is here!**

**PEPSI-SAXS and PEPSI-SANS**

# Assessing SAS-data model fits – Methods abound!

## Benchmarking predictive methods for small-angle X-ray scattering from atomic coordinates of proteins using maximum likelihood consensus data

Jill Trewhella,[a]* Patrice Vachette[b]* and Andreas Haahr Larsen[c]

[a]School of Life and Environmental Sciences, University of Sydney, NSW 2006, Australia, [b]Institute for Integrative Biology of the Cell (I2BC), Université Paris-Saclay, CEA, CNRS, Gif-sur-Yvette, Paris 91198, France, and [c]Department of Neuroscience, University of Copenhagen, Blegdamsvej 3, 2200 Copenhagen, Denmark. *Correspondence e-mail: jill.trewhella@sydney.edu.au, patrice.vachette@i2bc.paris-saclay.fr

EMBL

# Workhorse in ATSAS: CRYSOL (for SAXS)

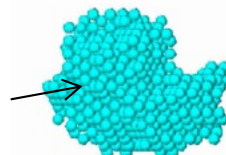Convert the atomic coordinates of a model into a convenient mathematical expression for fitting or modelling.

Calculate the envelope function from the centre of the macromolecule from a common/coincident grid origin.

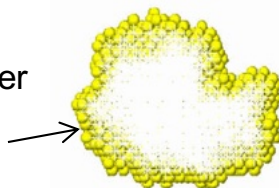**Take into account the atomic scattering, the excluded volume and hydration shell scattering.**

SLD of the solvent

$$I(s) = \left\langle |A(s)|^2 \right\rangle_\Omega = \left\langle |A_a(s) - \rho_s A_s(s) + \delta\rho_b A_b(s)|^2 \right\rangle_\Omega$$

Electrons (nuclei) are 'points'

Excluded volume with the SLD of the solvent.

SLD of the hydration layer is slightly different to bulk.

- $A_a(s)$: atomic scattering amplitudes in vacuum
- $A_s(s)$: scattering amplitudes from the excluded volume
- $A_b(s)$: scattering amplitudes from the hydration shell

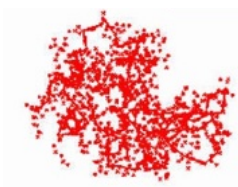**CRYSOL (X-rays):** Svergun et al. (1995). *J. Appl. Cryst.* **28**, 768
**CRYSON (neutrons):** Svergun et al. (1998) *P.N.A.S. USA*, **95**, 2267

EMBL

# Workhorse in ATSAS: CRYSOL

$$I(s) = \left\langle \left| A(s) \right|^2 \right\rangle_\Omega = \left\langle \left| A_a(s) - \rho_s A_s(s) + \delta\rho_b A_b(s) \right|^2 \right\rangle_\Omega$$

- Either fit the experimental data by varying the density of the hydration layer $\delta\rho$ (affects the third term) and the total excluded volume (affects the second term).
- Or predict the scattering from the atomic structure using default parameters (theoretical excluded volume and bound solvent density of $1.1 g/cm^3$).
- Provide output files (scattering amplitudes) for rigid body refinement routines.
- Compute particle envelope function $F(\omega)$

$$A(s) = \sum_{i=1}^{N} b_i \, e^{i\vec{s} \cdot \vec{r}}$$

Spherical wave bit

'Scattering factor': relates to the atomic cross section, i.e., scattering length, or probability of an atom to scatter for every atom in the sample.

The 'spherical wave bit' can be mathematically expressed in terms of a summed set of independent **spherical harmonics** (as a multipole expansion):

$$I(s) = 2\pi^2 \sum_{l=0}^{\infty} \sum_{m=-l}^{l} \left| A_{lm}(s) \right|^2$$

In 1970, Stuhrmann showed that the information content of a SAXS profile can be conveniently described in terms of a sum of spherical harmonic functions.

EMBL

# Spherical Harmonics



Essentially given a set of atomic coordinates in 3-dimensions (i.e., x, y, z coordinates), and knowing the identity of each atom at that coordinate (i.e., the atomic form factor), as well as the atomic volumes and scattering length densities, we can calculate the scattering amplitudes from the entire structure. As a result we can calculate the scattering intensities (i.e., the square of the scattering amplitudes.)

$$F(\omega) \cong F_L(\omega) = \sum_{l=0}^{L} \sum_{m=-l}^{l} f_{lm} \cdot Y_{lm}(\omega)$$

$$\rho(\mathbf{r}) = \begin{cases} 1, & 0 \le r \le F(\omega) \\ 0, & r > F(\omega) \end{cases}$$

$A_{00}(s)$    $A_{11}(s)$

$= \quad + \quad + \quad + \quad + \quad$ etc …

$A_{20}(s)$     $A_{22}(s)$

$I(s) = \langle I(s) \rangle$ = the Fourier transform of $\rho(r)$ squared i.e., $\langle (F\rho(r))^2 \rangle$

EMBL

# How many spherical harmonics to use in CRYSOL?

If you use the first harmonic only, i.e., zeroth-order, then the calculated intensities from the model will be a sphere. This is okay only if you want to describe the overall SIZE of the object, i.e., at the very lowest of angles in the Guinier region of the scattering profile. The zeroth-order harmonic dominates the very lowest angles of a calculated scattering profile!

If you use two harmonics, you will introduce an additional 'shape feature' into the calculated scattering intensities across *s*...but the resulting shape will probably still look like a sphere..with a couple of very low humps.

If you continue to increase the number of harmonics, you introduce additional shape features across *s*. However, the more harmonics you introduce the less impact on the overall calculated scattering is observed at the low angles (i.e., in the SAXS regeime).

Typically 15-30 harmonics are used to describe size and the shape of the object. However, this depends on the CLASSIFICATION of an object. Clearly, if the object is an extended rod, you probably need additional spherical harmonics terms.

EMBL

# *I*(*s*) from a globular structure using different numbers of harmonics

# *I*(*s*) from an extended structure using different numbers of harmonics



Lower order capture
Guinier

…yes this protein is real. It is
called SASG – SASBDB search it.

Higher-orders required to
describe the anisotropic
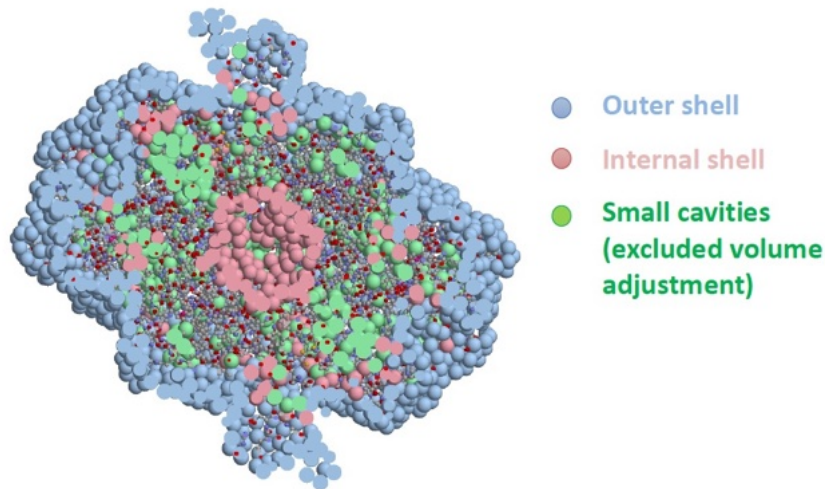structure! Computationally
expensive!

30 and 50

10

5

1

0

# Centre your atomic models!

- THE MODEL SCATTEING AMPLITUDES (and therefore the resulting intensities) MUST BE CALCULATED FROM THE ORIGIN, i.e., the models must be centred, otherwise you loose low-order harmonic contributions.

ATSAS tool: *alpraxin*

EMBL

# For macromolecules with cavities and holes – explicit hydration water using CRYSOL3

- Hydration shell representation as envelope function (CRYSOL – implicit solvent layer) or dummy solvent beads, i.e., explicit solvent layer (CRYSOL 3).
- Explicit solvent modelling is important for internal cavities!



- Outer shell
- Internal shell
- Small cavities (excluded volume adjustment)

Especially important for ring-shaped, hollow sphere, very small (less than 10 kDa) or very extended particles. Otherwise CRYSOL is fine.

EMBL

# Assessing data-model fits – $\chi^2$

…knowing what model does NOT fit the data can be as valuable as knowing what model(s) do fit the data!

CRYSOL fit to the SAXS data. The goodness of fit is described by the reduced $\chi^2$ discrepancy.

Calmodulin: X-ray crystal structure



PDB: 3CLN

$$\chi^2 = \frac{1}{N-1}\sum_j \left[ \frac{I_{exp}(s_j) - cI(s_j)}{\sigma(s_j)} \right]^2$$



$\chi^2 = 20.8$
The crystal structure does not fit the solution scattering whatsoever!

EMBL

# A note on $\chi^2$

$$\chi^2 = \frac{1}{N-1}\sum_j\left[\frac{I_{exp}(s_j) - cI(s_j)}{\sigma(s_j)}\right]^2$$

The errors on the scattering intensities need to be correctly specified, otherwise the test is, by default, INVALID. Errors follow Poisson counting statistics that limit to a gaussian distribution after many repetitions (for photon counting detectors).

If the errors are correctly specified and no significant (systematic) deviations are present between the experimental and modeled intensities, the value should lie in the range of approximately 0.9-1.1 depending on the number of points in the dataset (0.9-1.1 is typical for over-sampled SAXS data on modern detectors).

**Same intensities, same model, but different error estimates**

$\chi^2$ = 20.8
With correct errors

$\chi^2$ = 1.2
With incorrect errors



MBL

# Correlation Map: CorMap *P*

$$J = \begin{pmatrix} \vdots \\ I(q_k) \\ \vdots \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \ddots & & & \\ & \sigma(I(q_k))^2 & \cdots & \sigma(I(q_k),I(q_l)) \\ & \vdots & \ddots & \vdots \\ & \sigma(I(q_l),I(q_k)) & \cdots & \sigma(I(q_l))^2 \\ & & & & \ddots \end{pmatrix}$$

$$\sigma(I_{\exp}(q_k))^2 = \frac{1}{m-1}\sum_{i=1}^{m}(I_{\exp}(q_k)_i - \bar{I}_{\exp}(q_k))^2$$

On-diagonal variance.

$$\sigma(I_{\exp}(q_k),I_{\exp}(q_l)) =$$

$$\frac{1}{m-1}\sum_{i=1}^{m}(I_{\exp}(q_k)_i - \bar{I}_{\exp}(q_k))(I_{\exp}(q_l)_i - \bar{I}_{\exp}(q_l))$$

Off-diagonal co-variance between all point-
to-point $q_k$ and $q_l$.

Statistically significant, systematic differences between the modelled and experimental intensities

**P < 0.01**

View as a +/- 1 'map': random small patches = low probability of systematic differences (i.e., the pairwise comparison fits)!

**P > 0.01**



EMBL

# Error normalized residual plots

- Model fits to the data are also evaluated using normalized residual plots to help assess systematic model-fit deviations from the data in addition to over or under-estimation of the errors.

$$resudual = \frac{(I(s)experiment - cI(s)model)}{\sigma(I(s))}$$



$\chi^2$: 0.48

**Over-estimated errors**
CorMap P > 0.01

$\chi^2$: 0.01

**Severely over-estimated errors**
CorMap P < 0.01

$\chi^2$: 12.5

**Severely under-estimated errors**
CorMap P > 0.01

$\chi^2$: 1.0
CorMap P > 0.01

✔ *Correctly specified errors*

EMBL

# Lets do some atomistic model building!

SREFLEX
OLIGOMER
SASREF
BUNCH
CORAL
Ensemble Optimization Method (EOM)

EMBL

Modelling 3D-structures that fit SAXS data is perhaps the fundamental 'art' of small-angle scattering!

The major considerations to keep in mind when modelling SAS data are:

**There is often more than one model that fits the data equally well.**

## SAS data is inherently noisy.
## SAS data is inherently ambiguous.

Lets do the easy bit first: get the right sequence and the right PDB (or .cif) file(s).

- You should know the amino acid sequence of the protein (or polynucleotide, any other macromolecule, etc.) used for the SAS experiment. You should also know if the macromolecule binds metals, ligands, lipids, detergents, is glycosylated, etc.

- For proteins, use UNIPROT as a a fundamental resource to obtain the correct canonical sequence: www.uniport.org

- You should know what rigid-body (or bodies) you want to use for the modelling, i.e., the atomic coordinate PDB or .cif files.

- Extract the amino acid sequence from the PDB file.
- Align the atomic coordinate (.pdb/cif) amino acid sequence with the amino acid sequence of the **_EXACT_** protein used for the SAS measurement.
- Deal with missing side-chains in the atomic coordinate file (account for **_ALL OF THE MASS_**).

EMBL

**Amino acid sequence of protein used for SAS**

HMHHHHHHTRGSNNEEAICSLCDKKIRDRFVS
KVNGRCYHSSCLRCSTCKDELGATCFLREDS
MYCRAHFYKKFGTKCSSCNEGIVPDHVVRKA
SNHVYHVECFQCFICKRSLETGEEFYLIADDA
RLVCKDDYEQARDGGSGGHMGSGGGIGPLM
VQPATPHIDNTLGGPIDIQHF

↓

**Align the sequences using
Clustal Omega**

http://www.ebi.ac.uk/Tools/msa/clustalo/

↑

GSNNEEAICSLCDKKIRDRFVSKVNGRCYHSS
CLRCSTCKDELGATCFLREDSMYCRAHFYKK
FGTKCSSCNEGIVPDHVVRKASNHVYHVECF
QCFICKRSLETGEEFYLIADDARLVCKDDYEQ
ARDGGSGGHMGSGGGIGPLMVQPATPHIDNT
LGG
PIDIQHF

**Amino acid sequence of protein
from PDB or .cif file**

Atomistic model from PDB file (filename.pdb)



What is the amino acid sequence?

↓

ATSAS tool: *pdb2seq*

↓

This will save the sequence in the text file called
'filename.txt'

EMBL

Oops! Part of the sequence missing in the PDB file! This missing fragment will have to be built. Do not worry…ATSAS rigid-body modelling programs can deal with this!

**Amino acid sequence of protein used for SAS**

HMHHHHHHTRGSNNEEAICSLCDKKIRDRFVS
KVNGRCYHSSCLRCSTCKDELGATCFLREDS
MYCRAHFYKKFGTKCSSCNEGIVPDHVVRKA
SNHVYHVECFQCFICKRSLETGEEFYLIADDA
RLVCKDDYEQARDGGSGGHMGSGGGIGPLM
VQPATPHIDNTLGGPIDIQHF

↓

**Align the sequences using Clustal Omega**

http://www.ebi.ac.uk/Tools/msa/clustalo/

↑

GSNNEEAICSLCDKKIRDRFVSKVNGRCYHSS
CLRCSTCKDELGATCFLREDSMYCRAHFYKK
FGTKCSSCNEGIVPDHVVRKASNHVYHVECF
QCFICKRSLETGEEFYLIADDARLVCKDDYEQ
ARDGGSGGHMGSGGGIGPLMVQPATPHIDNT
LGG
PIDIQHF

**Amino acid sequence of protein from PDB or .cif file**

Atomistic model from PDB file (filename.pdb)



What is the amino acid sequence?

↓

ATSAS tool: *pdb2seq*

↓

This will save the sequence in the 'filename.txt'

Good to know and is useful, but this is old school!



Oops! Part of the sequence missing in the PDB file! This missing fragment will have to be built. Do not worry…ATSAS rigid-body modelling programs can deal with this!

EMBL

EMBL

# For proteins RNA, DNA, etc, just use AlphaFold3

https://golgi.sandbox.google.com/

# But my structure *almost* fits the data, can I just wiggle it a bit? - SREFLEX

Employs normal modes
pattern of motion on
domain-partitioned
structures.

Automated or manual
domain partitioning
possible.

Works with proteins!



Deciphering conformational transitions
of proteins by small angle X-ray
scattering and normal mode analysis

A. Panjkovich, D.I. Svergun (2016)
*Phys Chem Chem Phys.* 18, 5707-19

Used for spatial refinement of
models using small structural
adjustments.

Great for assessing whether slight
conformational movements are
required to fit SAXS data (e.g., from
crystal or AF-predicted structures).

Start structure

$\chi^2 = 2.6$

After normal mode

$\chi^2 = 1.2$

$I(s)$, a.u.

$s$, nm$^{-1}$

**ATSAS online version applied additional CONCORD refinement**

EMBL

Start structure

After normal mode

$\chi^2 = 2.6$

$\chi^2 = 1.2$

$I(s)$, a.u.

$s$, nm$^{-1}$

**ATSAS online version applied additional CONCORD refinement**

EMBL

# Combine SREFLEX with Multi-FoXS

Ask SREFLEX to output the normal mode models. Generate a normal mode pool. Do some basic scoring.

Get Multi-FoXS to fit the resulting NMA ensemble

Initial model: Manually define the rigid bodies in SREFLEX

Why FoXS? Well…I think it handles the somewhat complicated 'rough' glycosylated surface a bit better.

EMBL

# Scattering from mixtures

- Possible to obtain the volume fraction contribution to the total scattering profile of individual components of mixtures.

$$I(s) = \sum_k v_k I_k(s)$$

EMBL

# ATSAS program: OLIGOMER

Monomer model is too small: does not fit Guinier region!

Monomer fit: $\chi^2$ = 2.6
$R_g$ = 2.8 nm

Dimer model is too big: does not fit Guinier region!

Dimer fit: $\chi^2$ = 68!
$R_g$ = 3.9 nm

Experimental $R_g$ = 3.05

Oligomer fit: $\chi^2$ = 1.3
90% monomer; 10% Dimer

EMBL

# foXS combined with Multi-foXS!



Can upload a zip file with multiple structures

Assess the individual model fits, then also pass the models to Multi-foXS for oligomeric analysis

# Structure still does not fit – try some rigid body modelling

- The structures of two (or more) subunits in reference positions are known.

- Arbitrary complex can be constructed by moving and rotating the subunits.

- This operation depends on three Euler rotation angles and three Cartesian shifts.



EMBL

# Structure does not fit – try some rigid body modelling

- The structures of
  are known.

- Arbitrary comple
  the subunits.

- This operation d
  Cartesian shifts.



**Shift: x, y, z**

B

B'    A

**Rotation:**
$\alpha, \beta, \gamma$

The partial amplitudes of a rotated and displaced subunit are expressed *via* the initial amplitudes, three Euler rotation angles and three Cartesian shifts):

$$A^{(i)}{}_{lm}(s) = A^{(i)}{}_{lm}(s) \{A_0^{(i)}{}_{lm}(s), \alpha^{(i)}, \beta^{(i)}, \gamma^{(i)}, x^{(i)}, y^{(i)}, z^{(i)}\}.$$

$$I(s) = 2\pi^2 \sum_{l=0}^{L} \sum_{m=-l}^{l} |\sum_{n} A^n{}_{lm}(s)|^2$$

For symmetric particles, there are fewer parameters and the calculations are faster

Svergun, D.I. (1991). *J. Appl. Cryst.* **24**, 485-492

A

C

EMBL

The target function:

$$E(\{X\}) = \chi^2[(I(s), I_{\exp}(s)] + \sum_i \alpha_i P_i$$

is minimized…basically $\chi^2$ plus penalties!

Penalties describe model-based restraints and/or introduce the available additional information from other methods: MX, NMR, EM, Alphafold etc).

A brute force (grid) search is applied if the number of free parameters is small.

Otherwise a Monte-Carlo based technique (e.g. simulated annealing) is employed to perform the minimization of E({X}).

EMBL

# A note on $\chi^2$

$$\chi^2 = \frac{1}{N-1} \sum_j \left[ \frac{I_{\exp}(s_j) - cI(s_j)}{\sigma(s_j)} \right]^2$$

$$E(\{X\}) = \chi^2[(I(s), I_{\exp}(s)] + \sum_i \alpha_i P_i$$

EMBL

Incorporate information from EM, crystallography, NMR, biochemistry (e.g., cross linking, Mass-spec), FRET and bioinformatics…and of course for proteins…Alphafold!



+

Default 'sensible' modelling restraints like:

- Minimise clashes.

- Maintain contacts.

- Don't shift too far from the origin!

- For dummy residues, make dihedral angles and Ramachandran geometry sensible.

- Do not inter-penetrate subunits (interconnectivity).

EMBL

# SASREF (for SAXS), SASREFcv (for SAXS and SANS)

- Each subunit is treated as an individual rigid body. Protein, DNA, RNA, etc.

- Assumes the atomistic models are **COMPLETE i.e., no missing fragments or mass!**

- Options to perform **MIXTURE** modelling (e.g., monomer-dimer; SASREFmx) or **CONTRAST VARIATION** (SAXS and SANS; SASREFcv).

- Start from arbitrary initial orientations of the subunits – at the grid origin.

- Simulated annealing is employed.

- Search of interconnected spatial arrangement of the subunits without clashes.

- Random movement/rotation at one SA step.

- Fitting the scattering data by minimizing the target function.

- Additional restraints may be applied.

Petoukhov, M. V., and Svergun, D. I. (2006). *Eur Biophys J.*, 35, 567-576

EMBL

# SASREF restraints

Subunit clashes or disconnected models are penalised!

Inter penetrating
subunits are
penalised.

Disconnected
models are
penalised.

EMBL

# SASREF inputs

**For SAXS:**

- Rigid body starting models – centred to an origin. Protein, DNA, RNA, etc.
- Scattering amplitude files of each rigid-body model calculated using **CRYSOL**.
- Contacts file (optional).
- Symmetry information.

Alphafold3 of course be used if no contact information is available.

Contact information can be exceptionally useful!



No contact information



*Single* contact

EMBL

# Other docking methods

# More complicated examples – Dealing with missing stuff, linkers, etc.



X-ray crystal structure

EMBL

# Missing stuff, linkers, etc

Add glycans (ATSAS tool *glycosylation*)

Missing N-terminus – Dummy residue addition

Introduce movement between domains in the connecting linker

Missing C-terminus – Dummy residue addition

X-ray crystal structure

EMBL

# Missing stuff, linkers, etc



Add glycans (ATSAS tool *glycosylation*)

X-ray crystal structure

Missing N-terminus – Dummy residue addition

Introduce movement between domains in the connecting linker

Missing C-terminus – Dummy residue addition

EMBL

# BUNCH – will optimize domain and dummy-amino acid positions

- For SAXS only!

- Single residue **polypeptide chain only, i.e., 'protein domains'!**
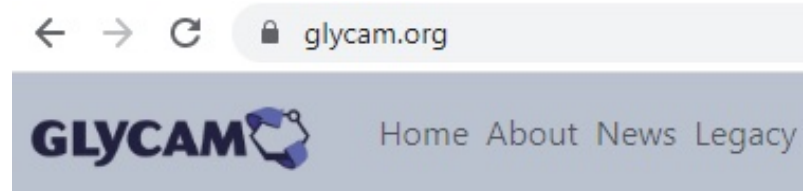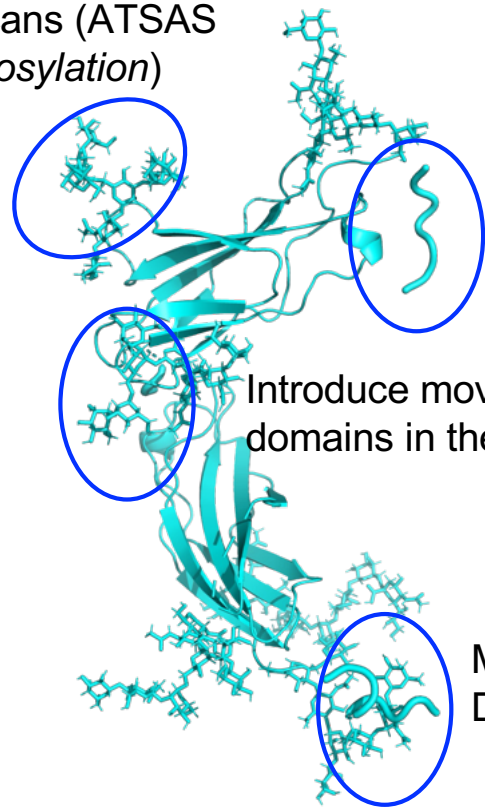
- With or without symmetry – **multiple curves allowed**, e.g., domain truncation mutants.

  - *Models missing linkers and mass as a set of dummy residues*.

- A two-step procedure.

    - pre_bunch
      - bunch

- Requires the domain PDB files and the EXACT protein sequence along with the SAXS data and scattering amplitudes calculated by CRYSOL.

EMBL

# BUNCH – will optimize domain and dummy-amino acid positions

- For

- Sing

- With
  trun

  - *M*

- A tw

- Rec
  (along with the SAXS data and scattering amplitudes calculated by
  CRYSC



Absence of
steric clashes

Neighbors
distribution
along the
sequence

Bond angles &
dihedrals distribution

Loop compactness
may also be required $Rg_{id} = 3\sqrt[3]{n_l}$

Petoukhov M.V., Svergun, D.I. (2005). *Biophys. J.* **89**, 1237-1250

EMBL

# ATSAS - CORAL

- SASREF – is good for modelling whole/complete complexes against SAXS data.

- BUNCH – is good form modelling single polypeptide chains with missing fragments, or linkers connecting modules/domains against SAXS data

  - *CORAL* combines both concepts into one!

- CORAL – Protein, DNA, RNA, glycosylated systems and complexes...all are possible!

- Known subunit interfaces can be preserved by grouping subunits together.

- CORAL is also a great deal faster than BUNCH (CORAL can be used to model single polypeptide chains as well, and it is much faster!).

- ...***SAXS only!!***

EMBL

- CORAL requires the SAXS data, domain/subunit atomic coordinate files along with the scattering amplitudes calculated by CRYSOL. A contact file is also possible!
- CORAL requires and additional .con file telling the program where to generate the linkers for each subunit:

NTER 6
KD_monomer1_1coral.pdb
LINK 10
KD_monomer2_1coral.pdb
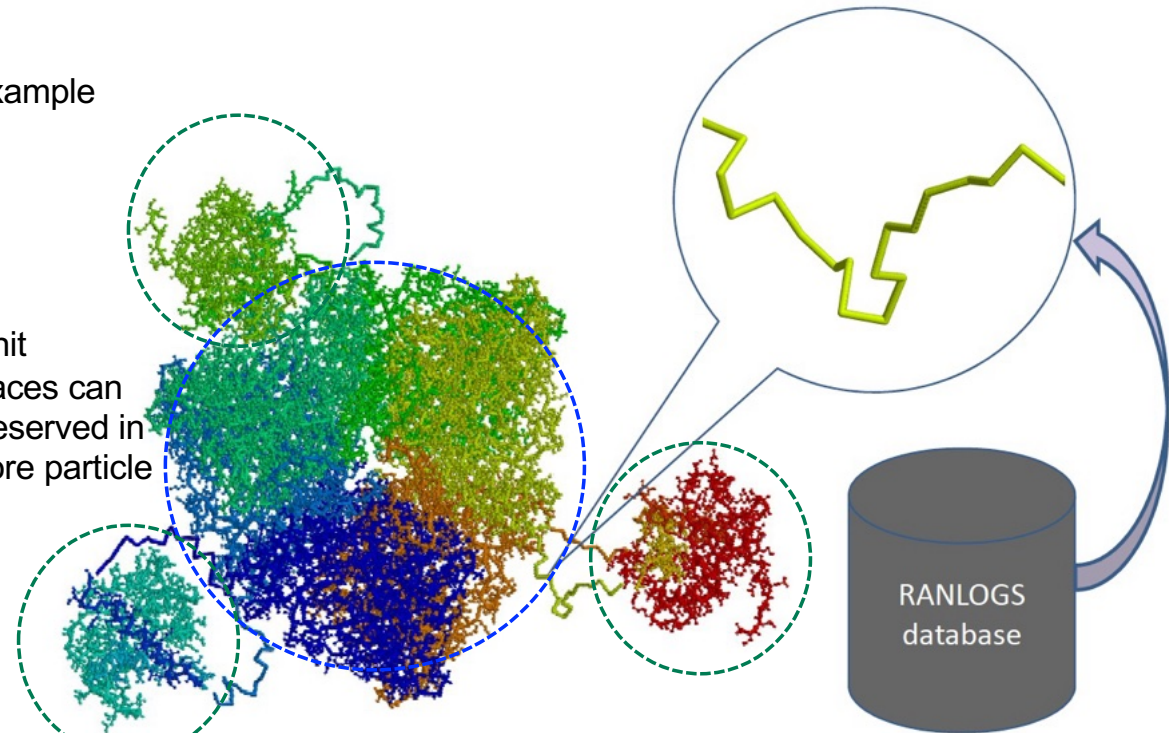CTER 10
NTER 4
KD_monomer2_coral.pdb
CTER 10
DNA.pdb

- At some point the program will ask:

`Pair of domains to group`

Where you can specify to preserve the spatial orientation between subunits. In the above, e.g., 3,4 and also 4,3 to preserve KD_monomer2_coral.pdb with DNA.pdb

For example

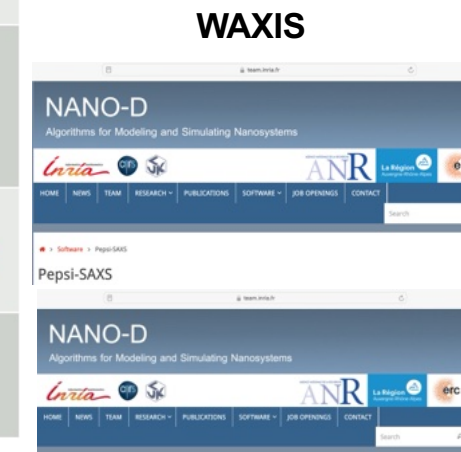Subunit interfaces can be preserved in the core particle



RANLOGS database

*CORAL*

Smaller extensions can be left to 'flop about' without obeying symmetry

EMBL

# Always check the final model fits using CRYSOL

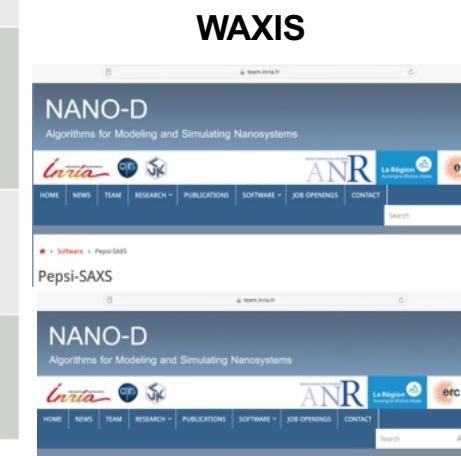| Approach | Modeling of the hydration layer | Representation of the molecule | References |
|---|---|---|---|
| CRYSOL | Implicit layer using an envelope function | All-atom | Svergun et al. *J. Appl. Cryst.* (1995) |
| AXES | Explicit water molecules using equilibrated water boxes | All-atom | Grishaev *et al. JACS* (2010) |
| FoXS | Implicit layer based on surface accessibility | All-atom or coarse-grained | Schneidman-Duhovny *et al. NAR* (2010) |
| HyPred | Explicit water molecules based on MD simulations | All-atom | Virtanen *et al. Biophys. J.* (2011) |
| AquaSAXS | Solvent-density map using the dipolar PB-Langevin approach | All-atom | Poitevin *et al. NAR* (2011) |

**WAXIS**

**PEPSI-SAXS and PEPSI-SANS** MBL

# Always check the final model fits using CRYSOL

| Approach | Modeling of the hydration layer | Representation of the molecule | References |
|---|---|---|---|
| CRYSOL | Implicit layer using an envelope function | All-atom | Svergun et al. *J. Appl. Cryst*. (1995) |
| AXES | | | *al.* |
| FoXS | | | *al.* |
| HyPred | Explicit water molecules based on MD simulations | All-atom | Virtanen *et al. Biophys. J*. (2011) |
| AquaSAXS | Solvent-density map using the dipolar PB-Langevin approach | All-atom | Poitevin *et al. NAR* (2011) |

You can use other programs as well, but be careful, **not all fitting programs can handle dummy-residues!**

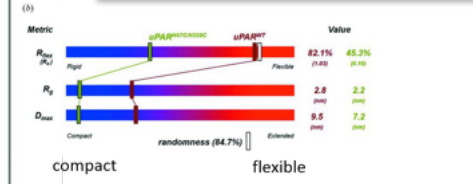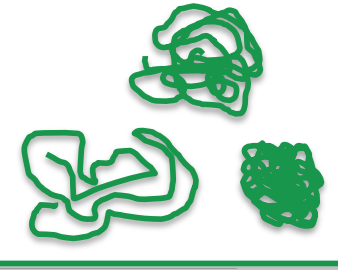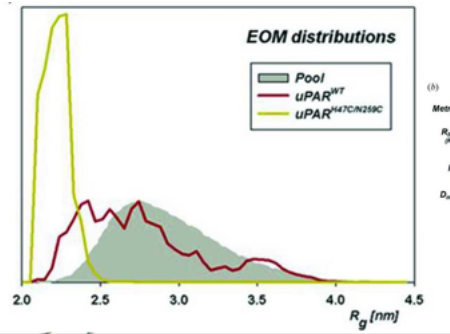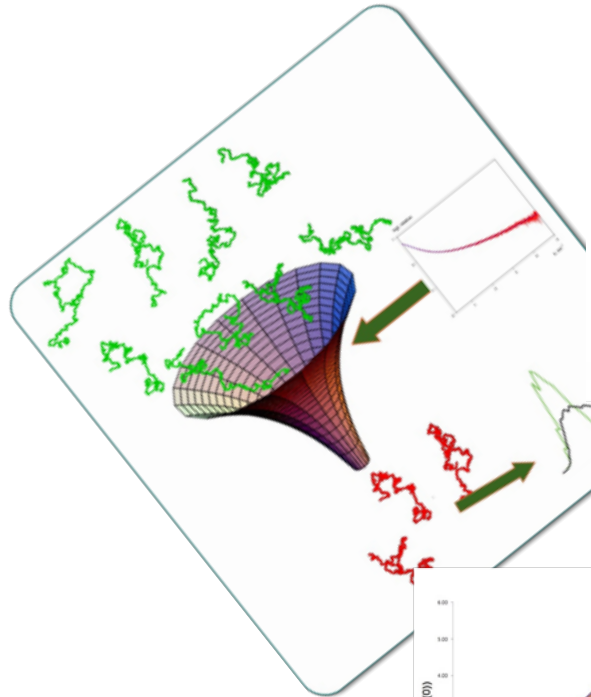**WAXIS**

**PEPSI-SAXS and PEPSI-SANS** MBL

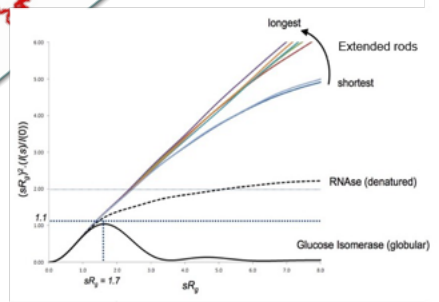# If there is one message, I want to get across today – always consider ambiguity!

- You **_must_** run your selected rigid body modelling routines at least 10 times and check for the spatial consistency of the models (spatial alignment using supcomb).

- *At the end of a BUNCH, SASREF or CORAL run check the fits with CRYSOL!*

- *Use Correlation Map to assess fits if you are unsure about your experimental errors!*

- *Error normalized residual plots are a great tool to visually assess systematic differences between modelled and experimental scattering intensities.*

- *…also apply common sense.*

- I usually do 20 modelling runs, check the individual model fits with CRYSOL (using 30 harmonics, minimum), then order the CRYSOL fits in terms of $\chi^2$ and CorMap $P$, then spatially align all models that fit the data to assess consistency.

EMBL

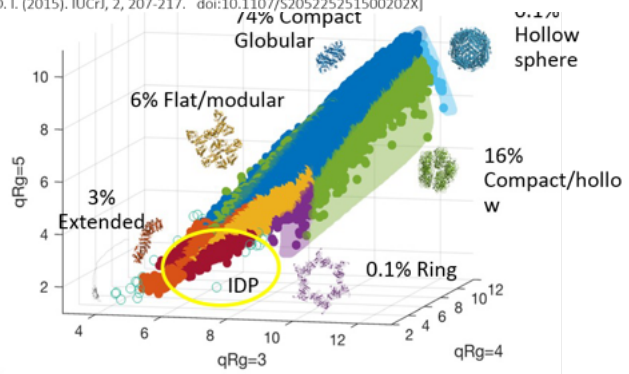# My structure is moving all over the shop!

Ensemble optimization method (EOM)



Characterization of the flexibility of uPARWT and the mutated uPARH47C-N259C using EOM 2.0. (a) Size distributions (Rg) of uPARWT and uPARH47C-N259C, providing only a qualitative assessment through direct comparison of the distributions of the selected ensembles and the pool. (b) The metrics Rflex and Rσ enable characterization of the flexibility quantitatively, with Rflex = ~82% and Rflex = ~45%, for uPARWT and uPARH47C-N259C, respectively, reflecting a significant change in compactness of the particle upon mutation (with a threshold of randomness of ~85% calculated from the pool). [Tria, G., Mertens, H. D. T., Kachala, M. and Svergun, D. I. (2015). IUCrJ, 2, 207-217. doi:10.1107/S205225251500202X]

Receveur-Bréchot & Durand (2012) *Current Protein and Peptide Science*, 13, 55-75.

Durand D, Vivès C, Cannella D, Pérez J, Pebay-Peyroula E, Vachette P, Fieschi F. (2010) *J Struct Biol.* 169: 45-53.
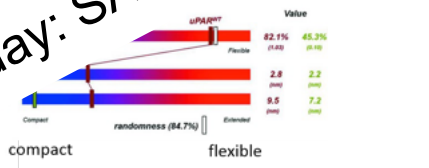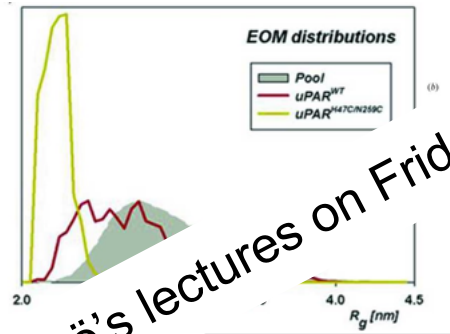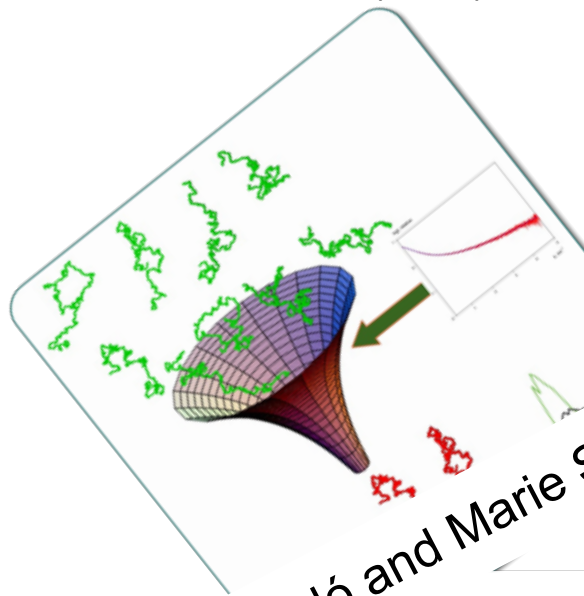
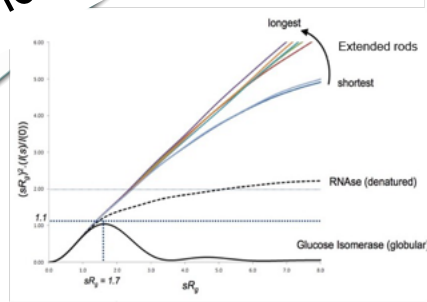Franke, Jeffries & Svergun (2018) *Biophys. J.* 114: 2485–2492

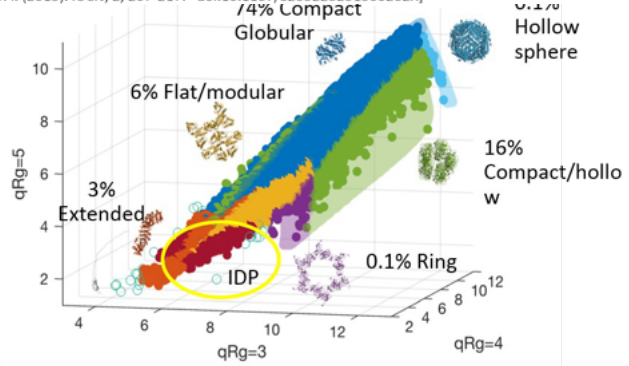# My structure is moving all over the shop!

Ensemble optimization method (EOM)



*Wait for Pau Bernadó and Marie Skepö's lectures on Friday: SAXS and flexibility*

SAXS and flexibility

**EOM distributions**

- Pool
- uPAR^WT
- uPAR^H47C/N259C

$R_g$ [nm]

compact    flexible

...of the flexibility of uPARWT and the mutated uPARH47C-N259C using EOM 2.0. (a) Size
...ns (Rg) of uPARWT and uPARH47C-N259C, providing only a qualitative assessment through
...comparison of the distributions of the selected ensembles and the pool. (b) The metrics Rflex and
...Rσ enable characterization of the flexibility quantitatively, with Rflex = ~82% and Rflex = ~45%, for
uPARWT and uPARH47C-N259C, respectively, reflecting a significant change in compactness of the particle
upon mutation (with a threshold of randomness of ~85% calculated from the pool). [Tria, G., Mertens, H.
D. T., Kachala, M. and Svergun, D. I. (2015). IUCrJ, 2, 207-217.  doi:10.1107/S205225251500202X]

Extended rods
longest
shortest

RNAse (denatured)
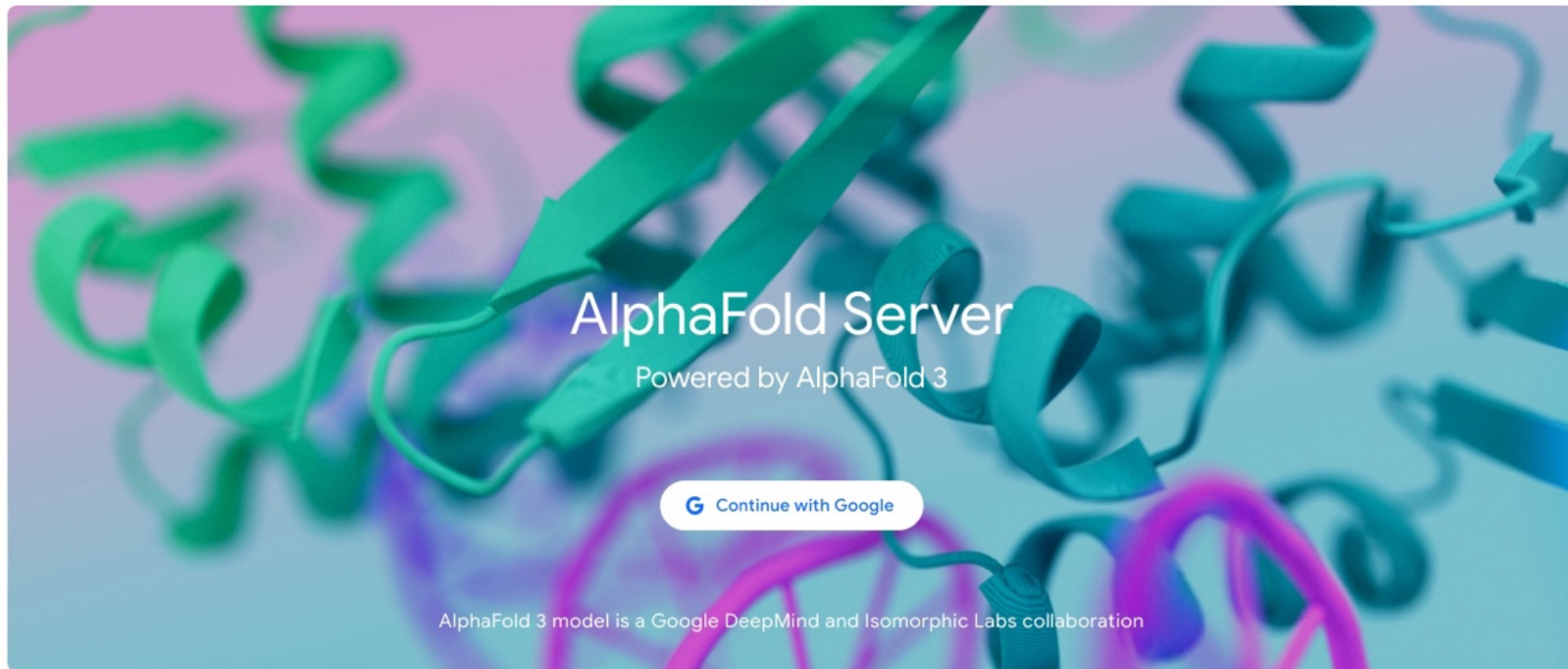
Glucose Isomerase (globular)

Receveur-Bréchot & Durand (2012) *Current Protein and
Peptide Science*, 13, 55-75.
Durand D, Vivès C, Cannella D, Pérez J, Pebay-Peyroula E,
Vachette P, Fieschi F. (2010) *J Struct Biol.* 169: 45-53.

74% Compact
Globular

0.1%
Hollow
sphere

6% Flat/modular

16%
Compact/hollow

3%
Extended

IDP    0.1% Ring

Franke, Jeffries & Svergun (2018) *Biophys. J.* 114: 2485–2492

EMBL

https://golgi.sandbox.google.com/

# Combine with Alphafold!

## AlphaFold-predicted protein structures and small-angle X-ray scattering: insights from an extended examination of selected data in the Small-Angle Scattering Biological Data Bank

Emre Brookes,[a]* Mattia Rocco,[b] Patrice Vachette[c] and Jill Trewhella[d]*
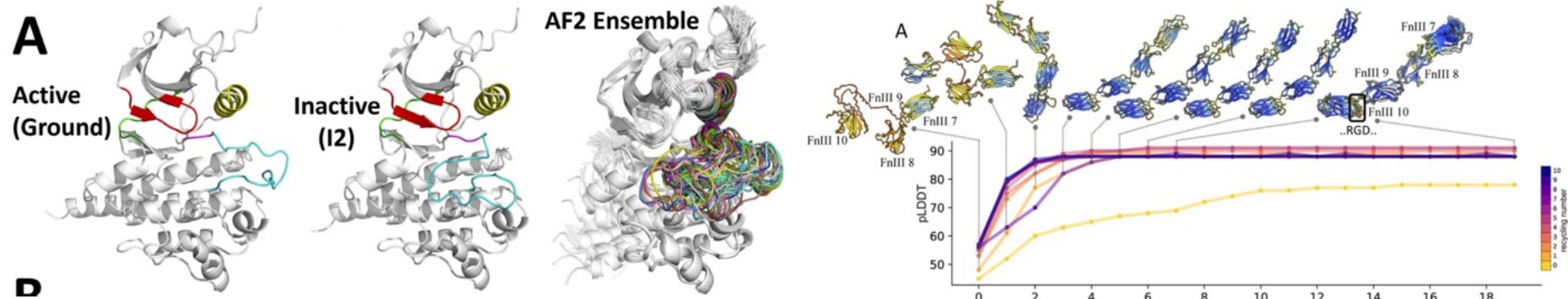
[a]Department of Chemistry and Biochemistry, University of Montana, 32 Campus Drive, Missoula, MT 59812, USA, [b]Proteomica e Spettrometria di Massa, IRCCS Ospedale Policlinico San Martino, Largo R. Benzi 10, Genova 16132, Italy, [c]Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), Gif-sur-Yvette 91198, France, and [d]School of Life and Environmental Sciences, The University of Sydney, NSW 2006, Australia. *Correspondence e-mail: emre.brookes@umontana.edu, jill.trewhella@sydney.edu.au

EMBL

# Combine with Alphafold!

AlphaFold-predicted protein st... and small-angle X-ray scattering: in... ... an extended examination of sele... in the Small-Angle Scattering Biolo... ...a Bank

Emre Bro... ...occo,[b] Patrice Vachette[c] and Jill Trewhella[d]*

...emistry and Biochemistry, University of Montana, 32 Campus Drive, Missoula, MT 59812, USA, ...e Spettrometria di Massa, IRCCS Ospedale Policlinico San Martino, Largo R. Benzi 10, Genova 16132, Italy, ...ersité Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), Gif-sur-Yvette 91198, France, and School of Life and Environmental Sciences, The University of Sydney, NSW 2006, Australia. *Correspondence e-mail: emre.brookes@umontana.edu, jill.trewhella@sydney.edu.au

*Wait for Emre Brookes and Rob Rambo lectures on Thursday: AI in SAXS*

EMBL

# Is it only a matter of time before AlphaFold can build ensemble models by itself?   …yes

**Thank you and goodbye!**

SAXS Team@EMBL

Dmytro Soloviov
Melissa Gräwert
Clement Blanchet
Aleksi Sutinen

Everyone involved in ATSAS over the years!