# SANS data analysis II: ATSAS and other methods

## for biological SAS (but only)

Wojtek Potrzebowski

# What SasView does not cover

- Ab initio modeling
- Rigid body modeling
- Shape and Conformational Polydispersity
- Efficient intensity calculation from PDB file
- Molecular Dynamics
- 3D particle electron densities from SAS data
- And much more smallangle.org/content/software
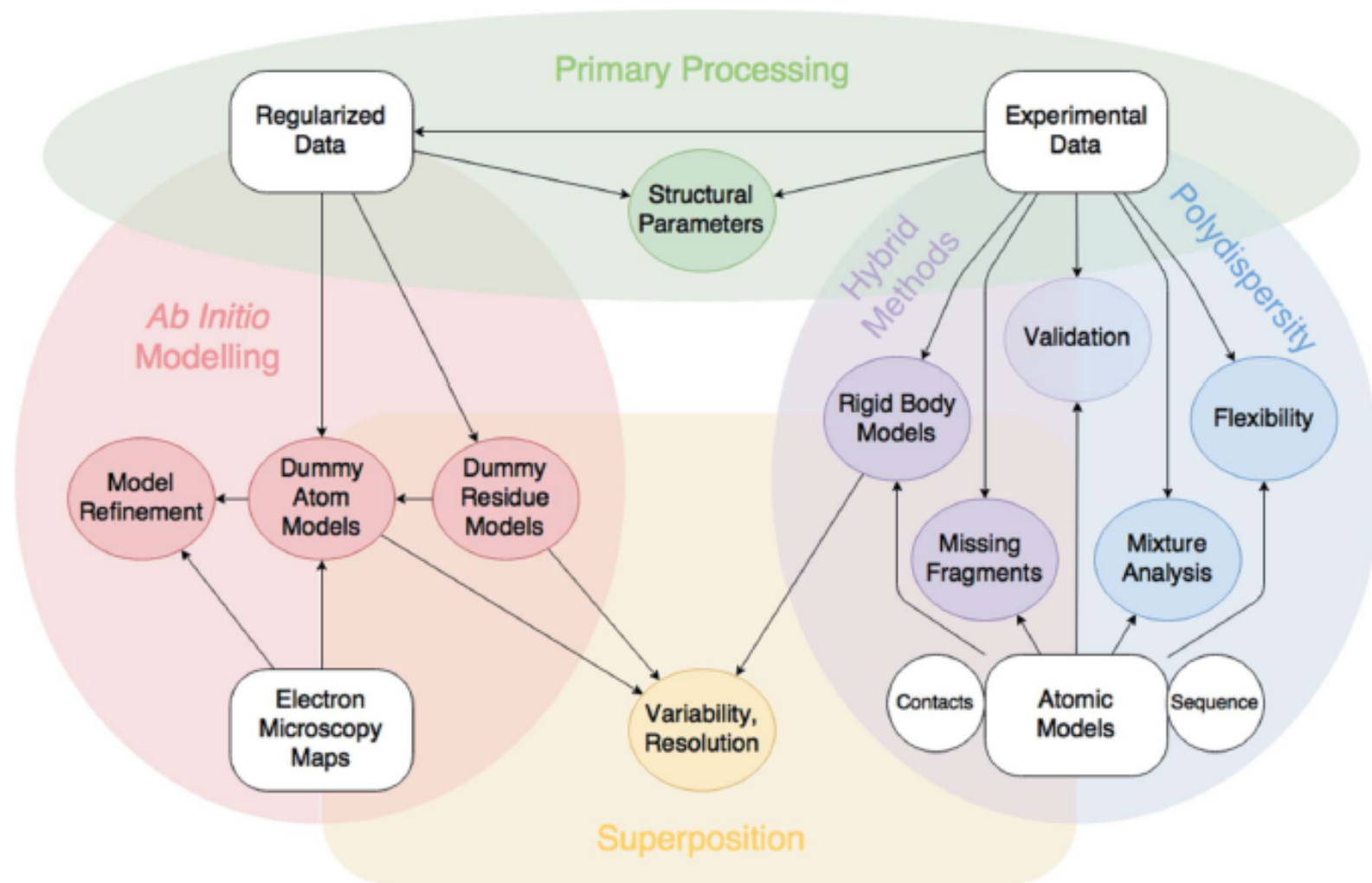
ATSAS
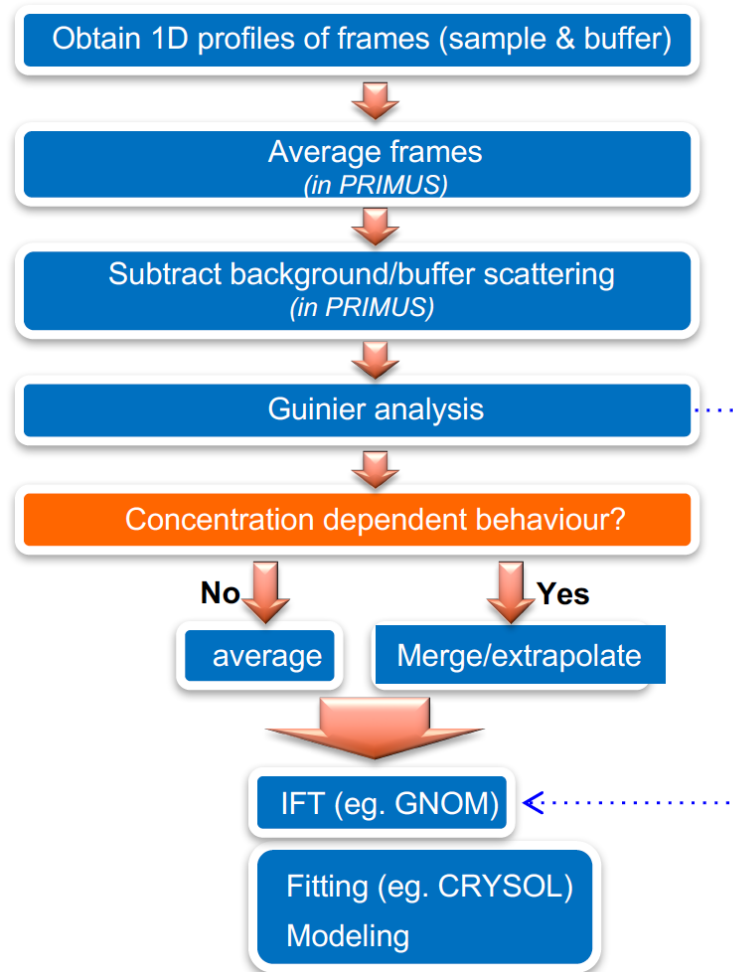
Sassie

DENSS

# 🄰 ATSAS software package 3.0

- Over 90 programs

- Operating systems:
  - Windows 8 and 10,
  - macOS 10.12 Sierra, 10.13 High Sierra and 10.14 Mojave,
  - Red Hat/CentOS 7 and 8,
  - Ubuntu 16 and 18,
  - Debian 9 and 10.

- Free for academic users:
  https://www.embl-hamburg.de/biosaxs/download.html

K. Manalastas-Cantos, P.V. Konarev, N.R. Hajizadeh, A.G. Kikhney, M.V. Petoukhov, D.S. Molodenskiy, A. Panjkovich, H.D.T. Mertens, A. Gruzinov, C. Borges, C.M. Jeffries, D.I. Svergun and D. Franke (2020) **ATSAS 3.0: Expanded functionality and new tools for small-angle scattering data analysis** *J. Appl. Cryst., submitted*
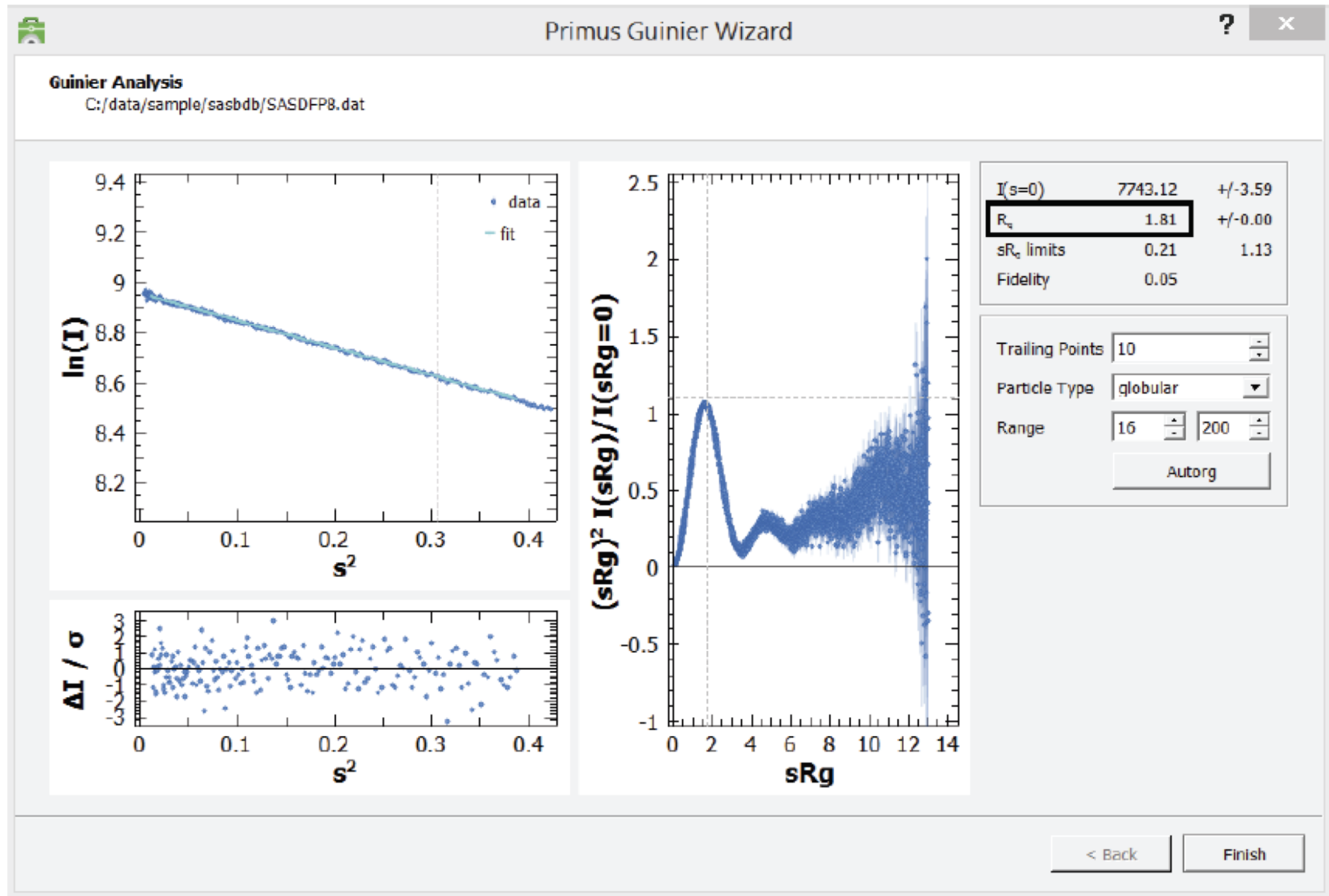
# Atsas software overview



Franke et al. J.Appl.Cryst. (2017). 50, 1212–1225

# Primary processing workflow



Adapted from Haydyn Mertens (EMBL-Hamburg), EMBO 2017

# Guinier approximation in PRIMUS



Adapted from Al Kikhney

# Molecular Weight in PRIMUS

# P(r) inversion
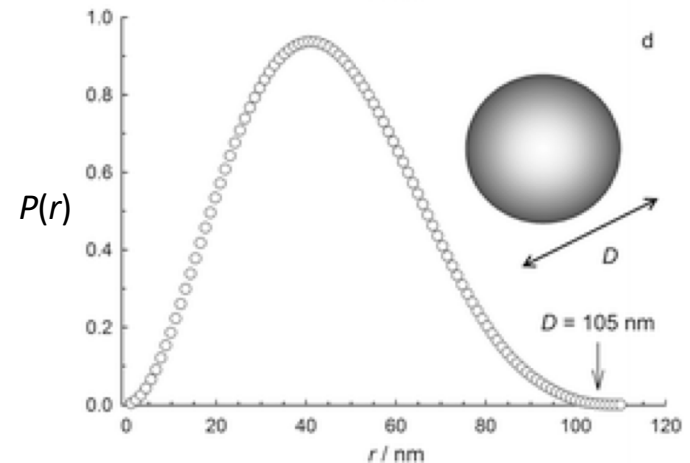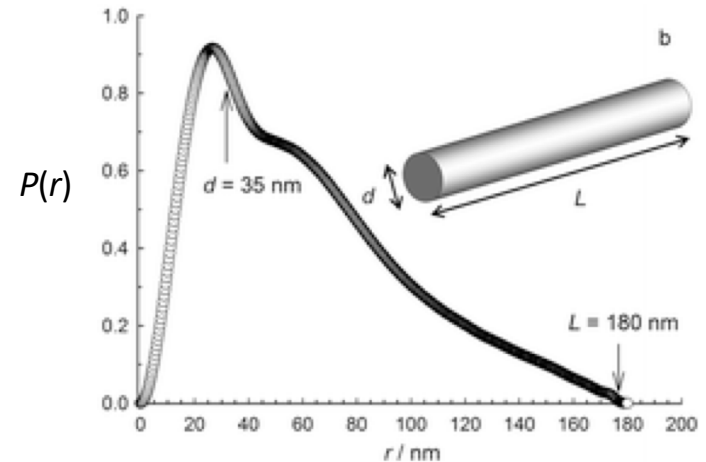
- Calculates a real-space pair-distance distribution function

$$I(q) = 4\pi \int_0^{D\max} p(r)\frac{\sin(qr)}{qr}.dr$$

**FT** ⇅ **FT⁻¹**

$$p(r) = \frac{r^2}{2\pi^2}\int_0^\infty q^2 I(q)\frac{\sin(qr)}{qr}.dq$$

- Calculated by Indirect Fourier Transform (Fourier transform of  noisy data).

- Popular methods: Glatter, Moore

- Maximum *d*, noise level, regularization constants have to be chosen

$P(r)$

$d = 35$ nm

$L = 180$ nm

$r$ / nm

$P(r)$

$D = 105$ nm

$r$ / nm

F. Grrohn, Soft Matter, 2010,6, 4296-4302

8

# P(r) calculation (SasView)

- Using IFT method (Moore, 1980)
- *P(r)* is set to be equal to an expansion of base functions of the type

$$P(r) = \sum_{n=0}^{N} c_n \Phi_n \qquad \Phi_{n(r)} = 2r\,sin\left(\frac{\pi n r}{D_{max}}\right)$$

- The coefficient $c_n$ of each base function in the expansion is found by performing a least square fit

$$\chi^2 = \frac{\sum_i (I_{meas}(Q_i) - I_{th}(Q_i))^2}{error^2} + Reg\_term$$

# P(r) calculation SasView

- *Number of terms*: the number of base functions in the P(r) expansion.

- *Regularization constant*: a multiplicative constant to set the size of the regularization term.

- *Maximum distance*: the maximum distance between any two points in the system.

# P(r) calculation (ATSAS)

**Indirect Fourier Transformation (IFT) of SAXS data**

- Solution is Indirect Fourier Transformation (IFT), (Glatter, 1977)

- Fit a function to the SAXS data and transform → *p(r)*

- Regularisation parameter (*α*) helps balance between the fit and the **FT**.

$$p(r) = \sum_{k=1}^{K} c_k \phi_k(s_i)$$
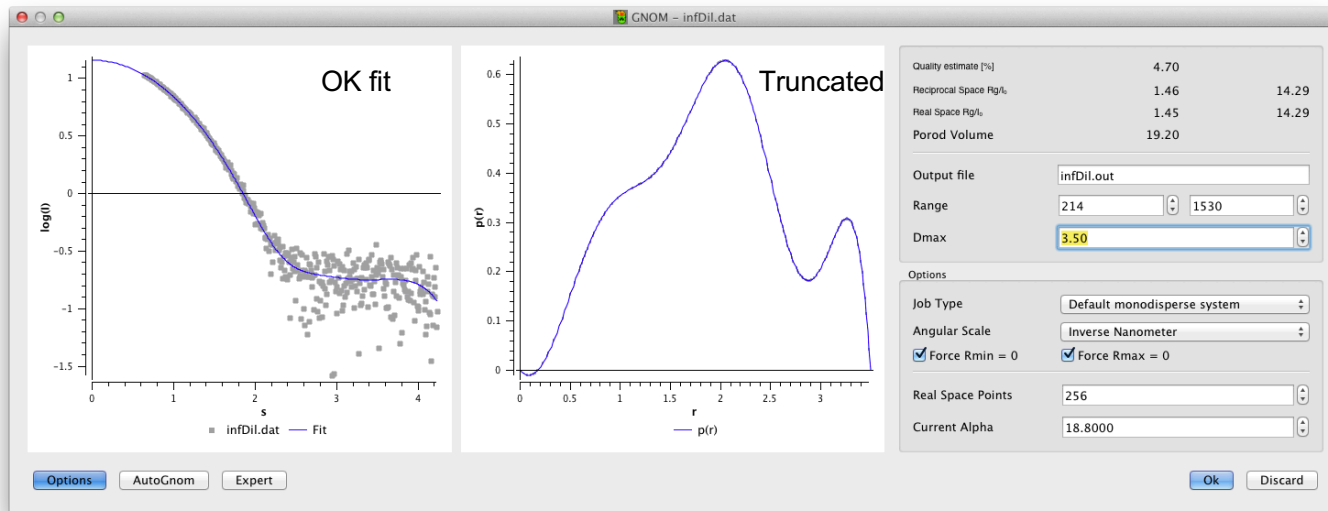
$$\Phi = \chi^2 + \alpha P(p)$$

$$P(p) = \int_0^{D_{max}} \left[ p' \right]^2 dr$$

$$\chi^2 = \frac{1}{N-1} \sum_{j=1}^{N} \left[ \frac{I_{exp}(s_j) - c I_{calc}(s_j)}{\sigma(s_j)} \right]^2$$

# P(r) calculation with GNOM

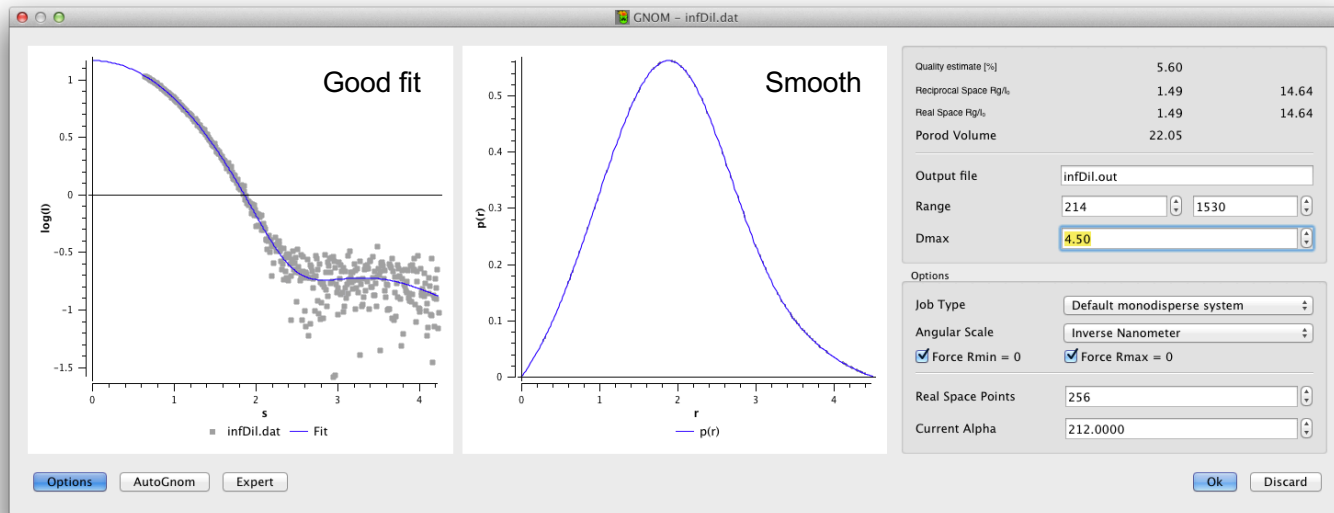**So, what is a good *p(r)*? How do I know a good solution?**
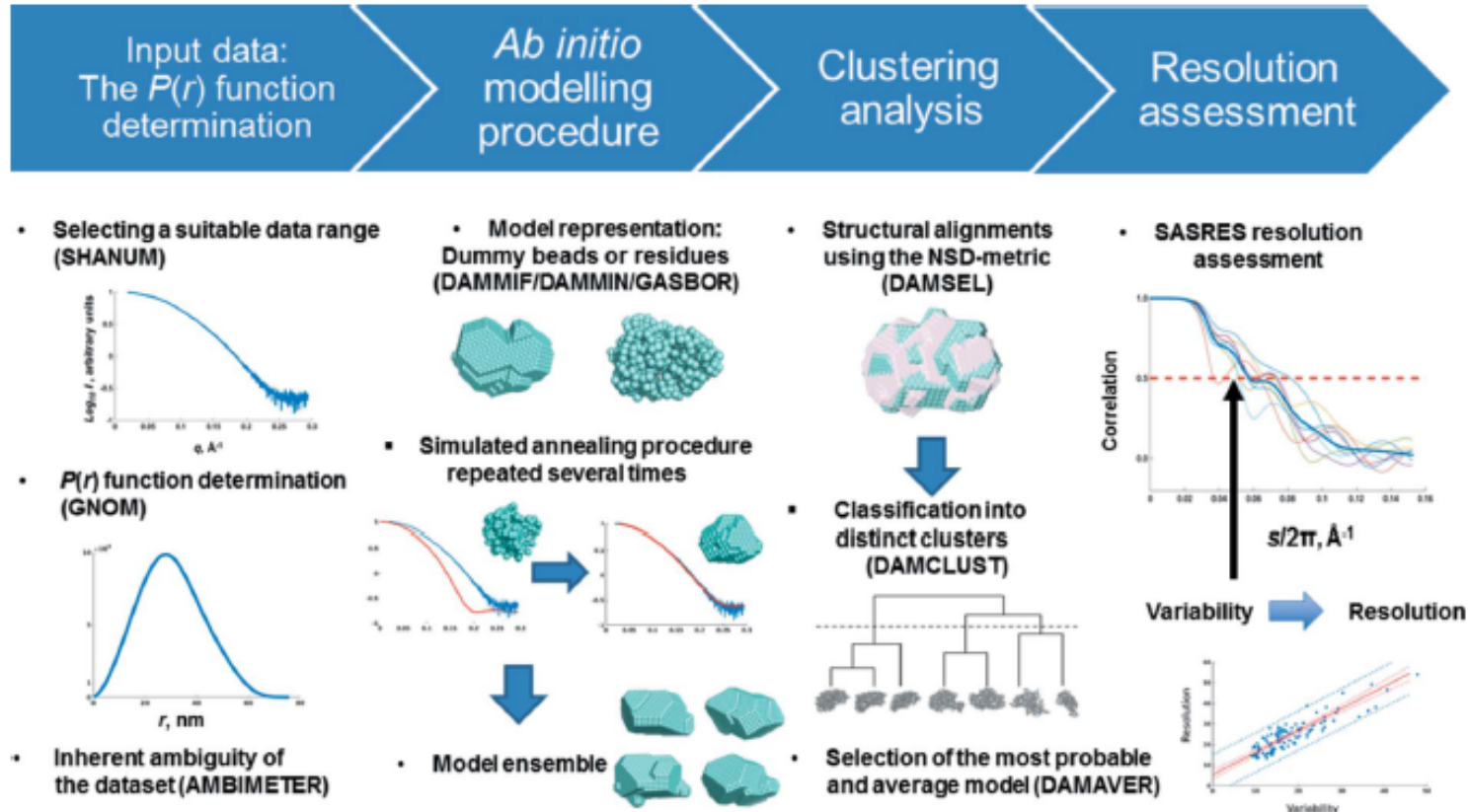
- *$D_{max}$* estimate (3.5 nm) poor solution – ***too small***

# P(r) calculation with GNOM

**So, what is a good *p(r)*? How do I know a good solution?**

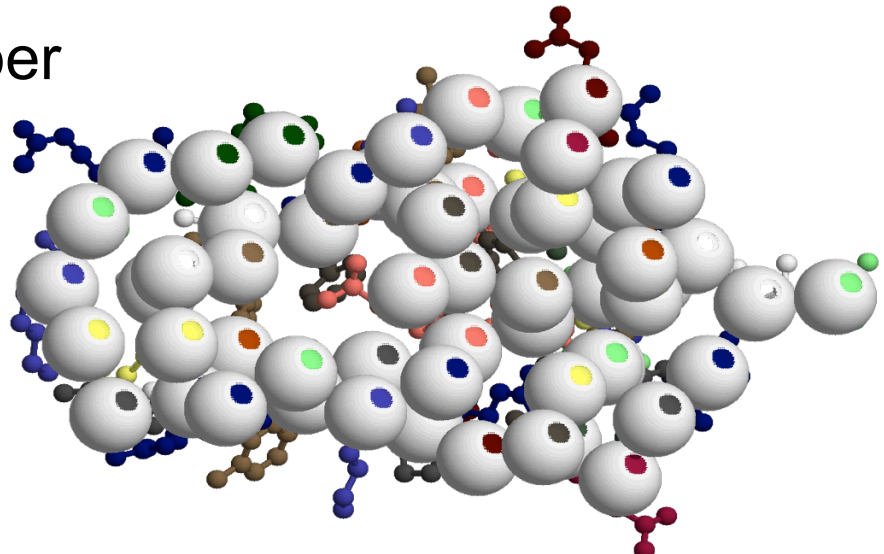- *$D_{max}$* estimate (4.5 nm) - ***good solution***
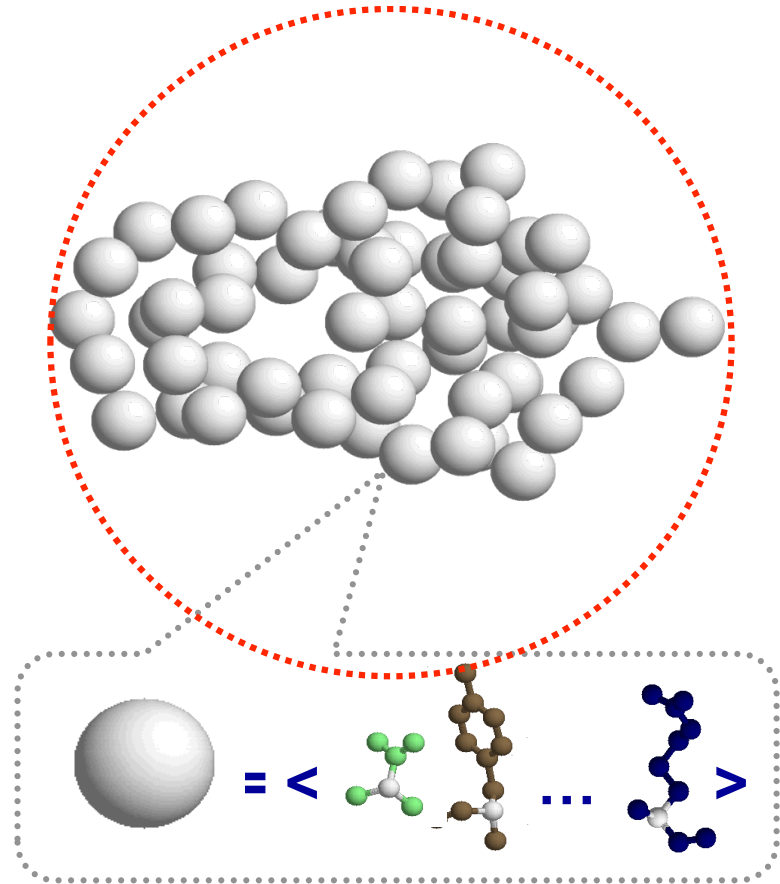
# Ab initio modeling overview

# Dummy residues

- Proteins typically consist of folded polypeptide chains composed of amino acid residues
- At a resolution of 0.5 nm each amino acid can be represented as one entity (dummy residue)
- In GASBOR a protein is represented by an ensemble of *K* dummy residues that are
  - Identical
  - Have no ordinal number
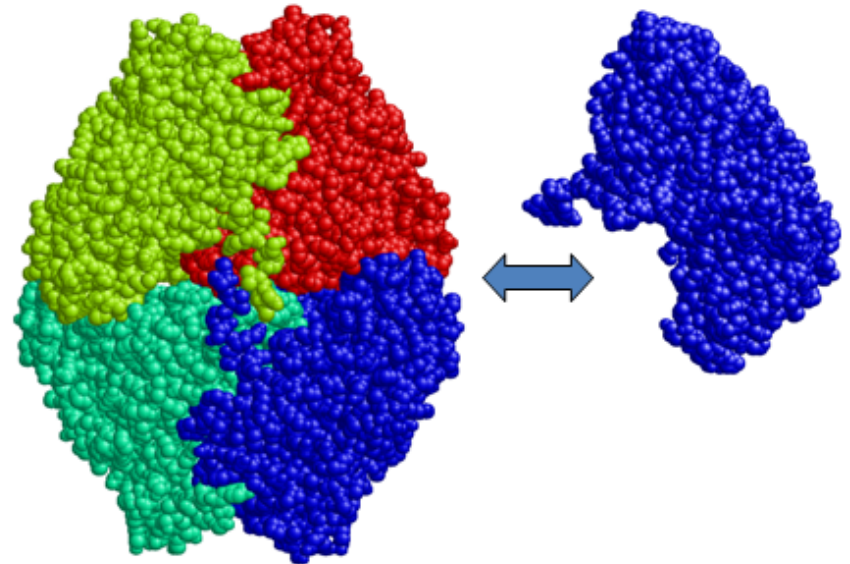  - For simplicity are centered at the C$\alpha$ positions

# Dummy residues

- GASBOR finds coordinates of $K$ dummy residues within its search volume (red)

- Scattering is computed using the Debye (1915) formula

- Requires polypeptide chain-compatible arrangement of dummy residues

# Dummy residues for mixture models

- GASBORMX extension to equilibrium mixtures

- Reconstructs the monomer and a symmetric multimer together

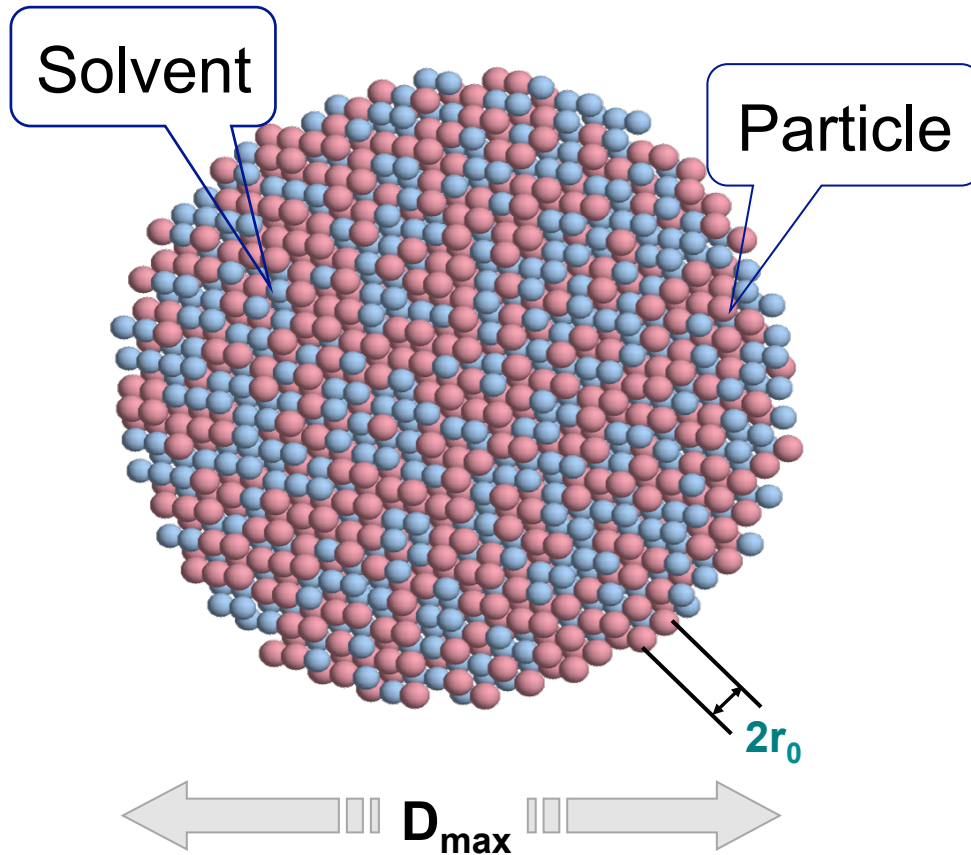- Interconnectivity is required for the monomer and the multimer

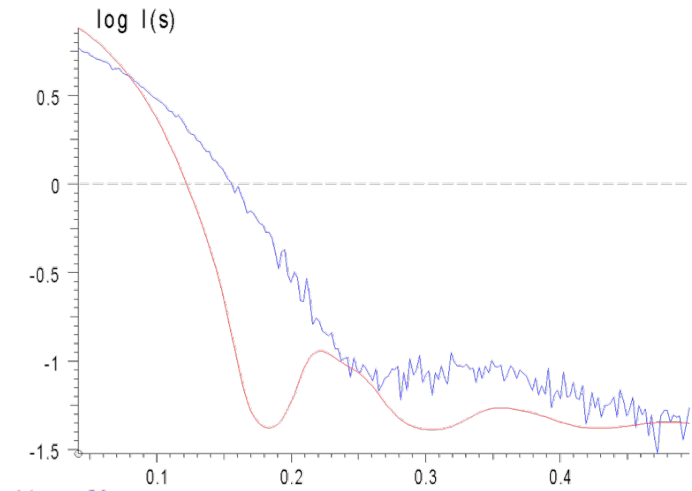# Single Phase Dummy Atom Models

Dummy atoms:

- Act as a placeholder for, but does not resemble, a real atom
- Occupy a known position in space
- Have a known scattering pattern
- May either contribute to solvent or particle
- Are also known as beads

# Single Phase Dummy Atom Models

A volume is filled by densely packed beads of radius $r_0 << D_{max}$



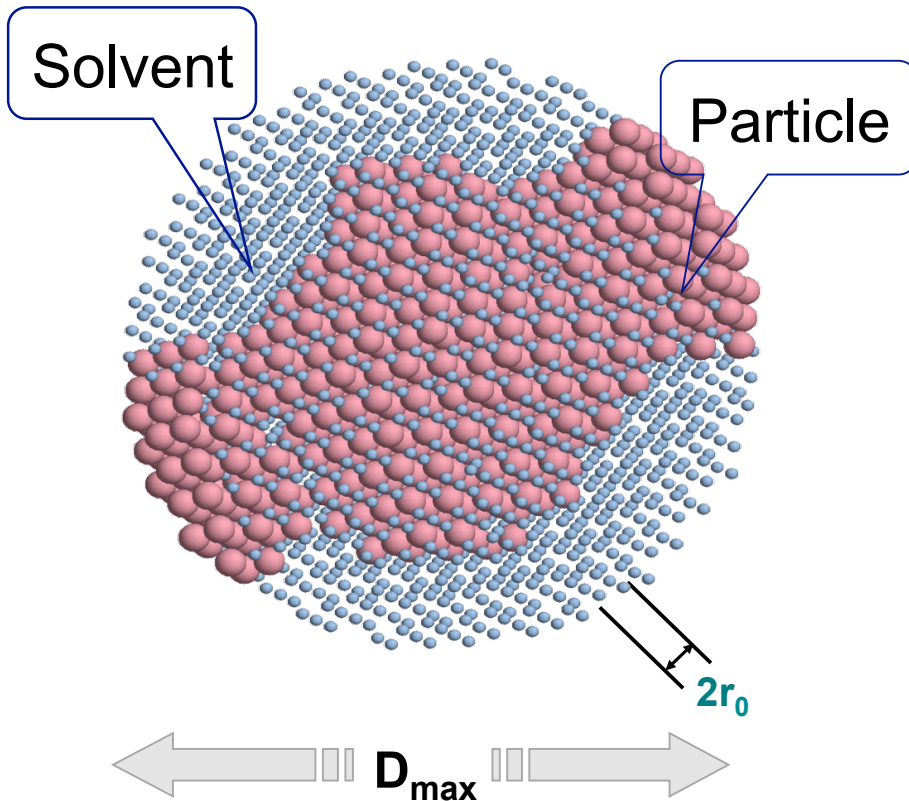Solvent

Particle

$2r_0$

$D_{max}$

Parametrization:
a binary vector,
0 if solvent, 1 if particle

# Single Phase Dummy Atom Models

A volume is filled by densely packed beads of radius $r_0 << D_{max}$



Solvent
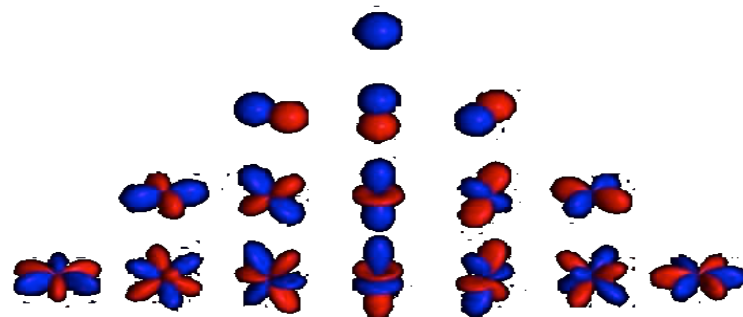
Particle

$2r_0$

$D_{max}$

Parametrization:
a binary vector,
0 if solvent, 1 if particle

# Single Phase Dummy Atom Models

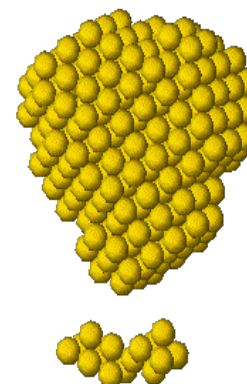- Scattering intensity is computed using spherical harmonics

- Penalty terms ensure compactness and connectivity

*compact*

*loose*
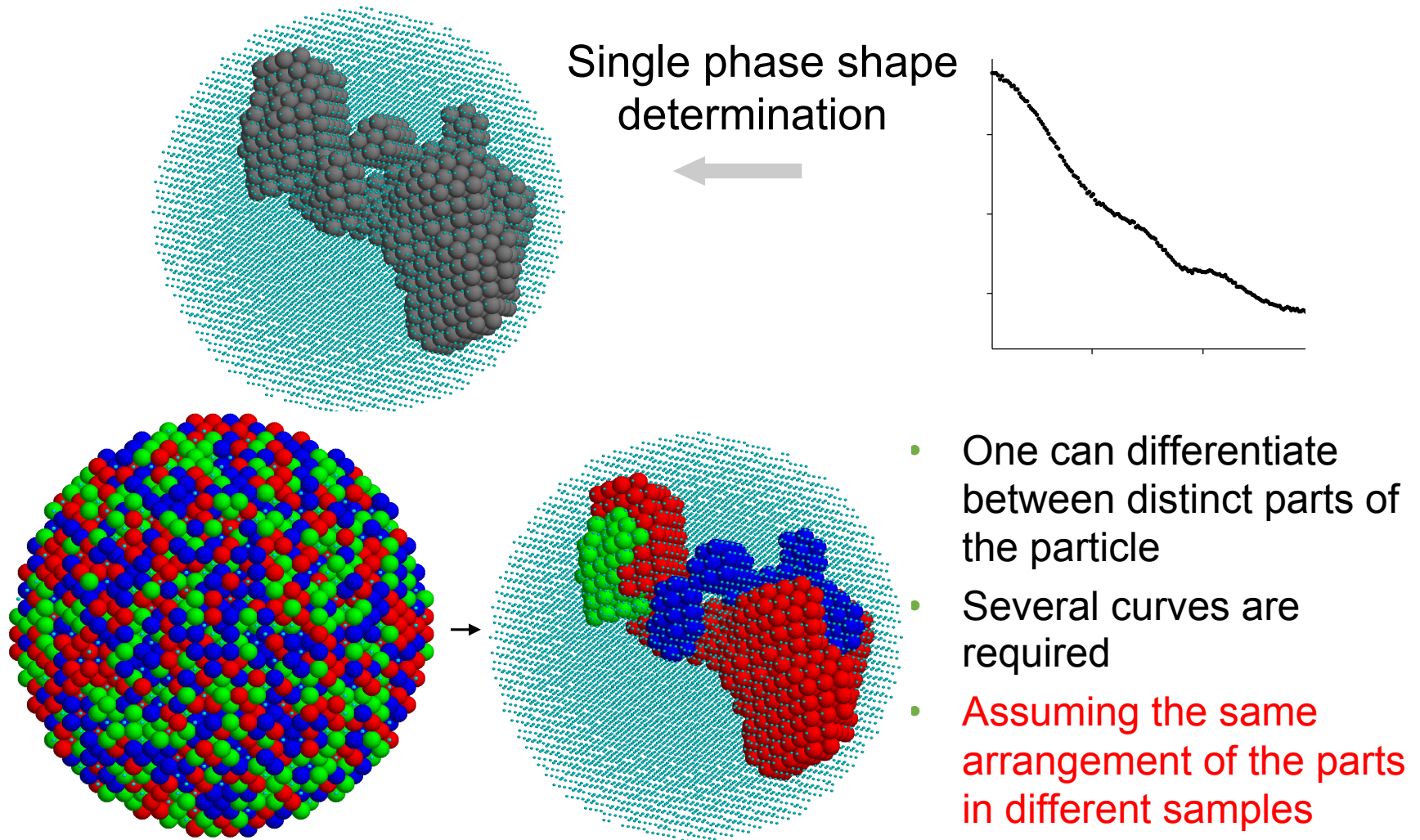
*disconnected*

# Multi Phase Dummy Atom Models



Single phase shape determination

- One can differentiate between distinct parts of the particle
- Several curves are required
- Assuming the same arrangement of the parts in different samples

Adapted from Daniel Franke (EMBL-Hamburg), EMBO 2017

# Dummy Atom Models

| | DAMMIN | DAMMIF | MONSA |
|---|---|---|---|
| Objects | any | any | any |
| Max # of phases | 1 | 1 | 4 |
| Angular range | lower part | lower part | lower part |
| Resolution | low | low | low |
| Search volume | fixed | growing | fixed |
| Constrains | Symmetry, Interconnectivity, Compactness | Symmetry, Interconnectivity, Compactness | Symmetry, Interconnectivity, Compactness |
| Performace | slow | fast | very slow |
| Limitations | | DAMMIN has better symmetry support | |

Warning: results are not atomic models, just a filled volume!
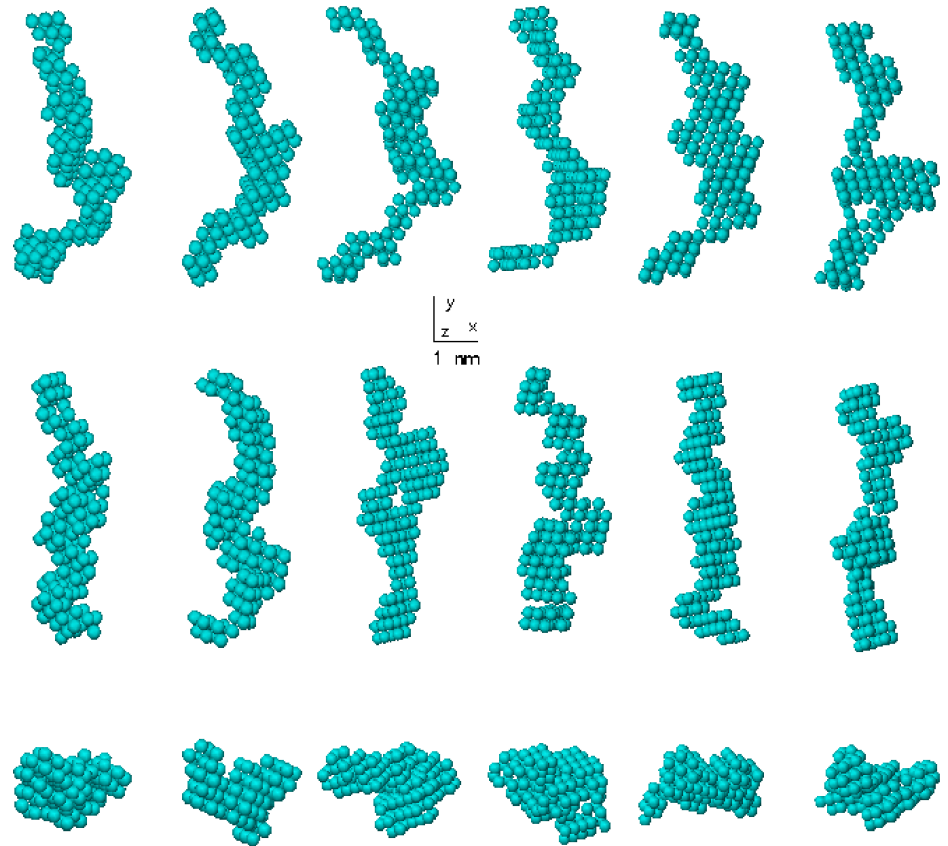
# Applicability to SAS data

| Program | SAXS | SANS |
|---|:---:|:---:|
| GASBOR/GASBORMX | ✓ | ✗ ** |
| DAMMIN/DAMMIF | ✓ | ✓ * |
| MONSA | ✓ | ✓ |

\* May be used if contrast is high and the particle is homogeneous
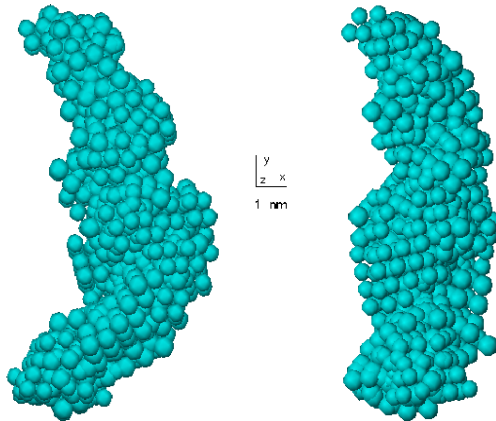** Dummy residue form factors are available for X-rays only
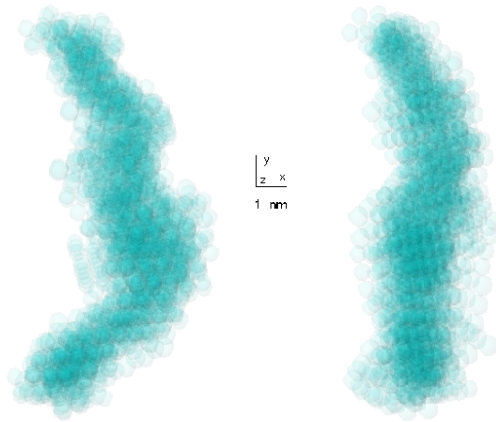
# Model post processing

- 5S RNA models
- A variety of DAMMIN models explains data equally well
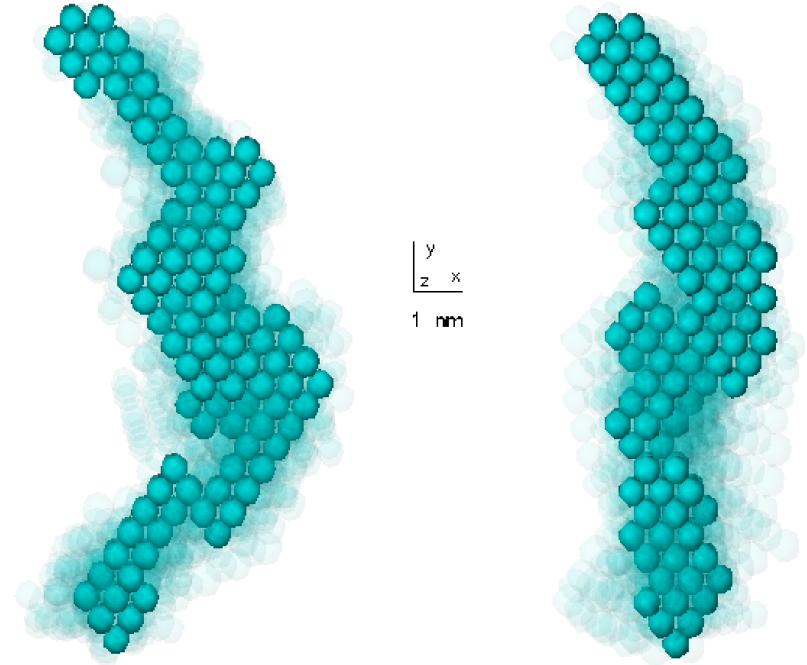
# Model post processing



5S RNA – Solution spread region



5S RNA – Most Populated Volume



5S RNA – Final Solution within the Spread Region

# Superposition with hi-res model



**SUPALM**

# Model Validity

- Validate your input data
- Check for
  - Aggregation
  - Noise at higher angles
- Keep in mind: it is easy to model noise


→ Garbage in, garbage out

# Fitting high-resolution structures to SAS data

# Scattering from macromolecule in solution

$$I(s) = \left\langle |A(s)|^2 \right\rangle_\Omega = \left\langle |A_a(s) - \rho_s A_s + \delta\rho_b A_b(s)|^2 \right\rangle_\Omega \qquad (2)$$
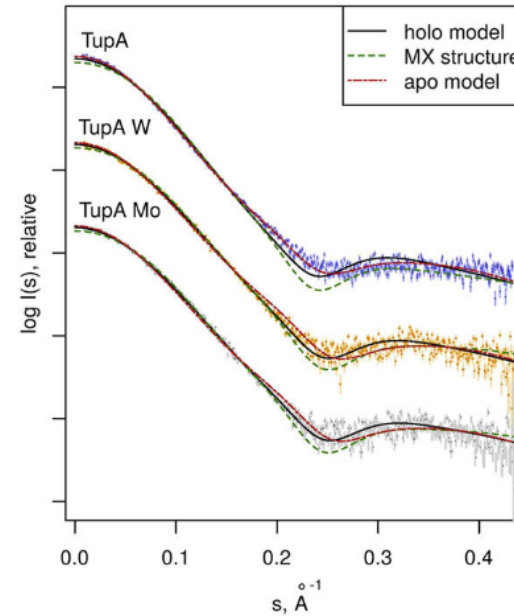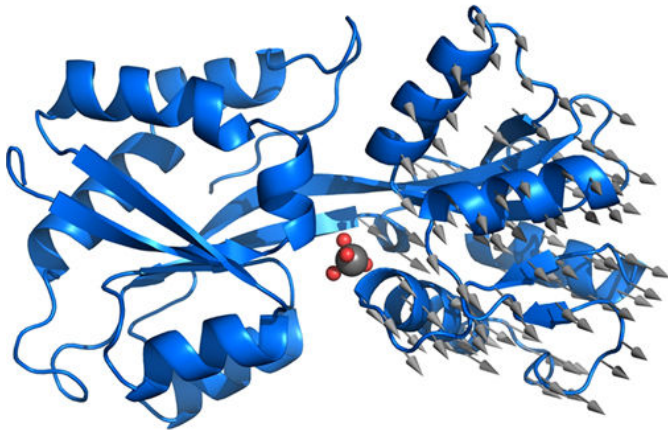
- $A_a(s)$: atomic scattering in vacuum
- $A_s(s)$: scattering from the excluded volume
- $A_b(s)$: scattering from the hydration shell

Programs:

- 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

# CRYSOL on PDB structure

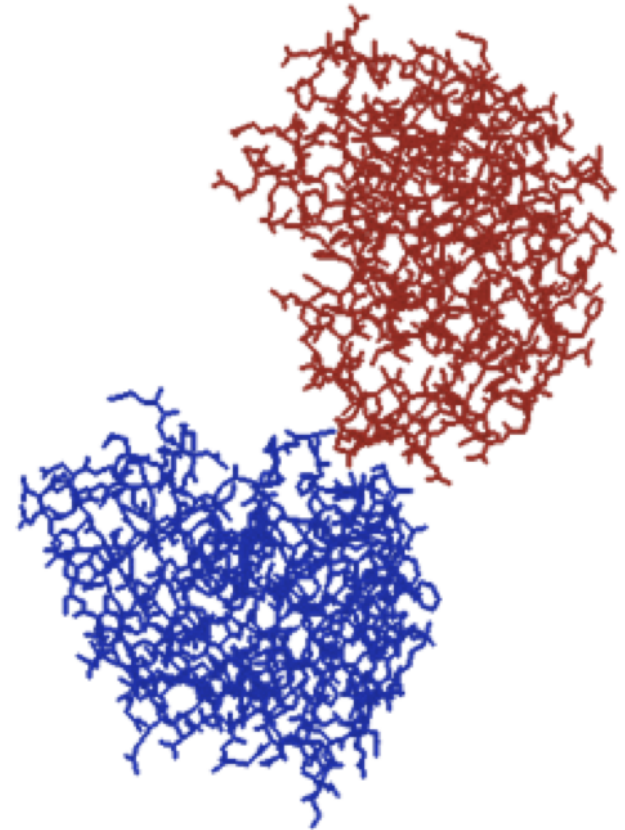How does the atomic model fit the solution scattering profile?



$$\chi^2 = \frac{1}{N} \sum_{i=1}^{N_p} \left( \frac{I_e(s_i) - cI(s_i)}{\sigma(s_i)} \right)^2 \tag{7}$$
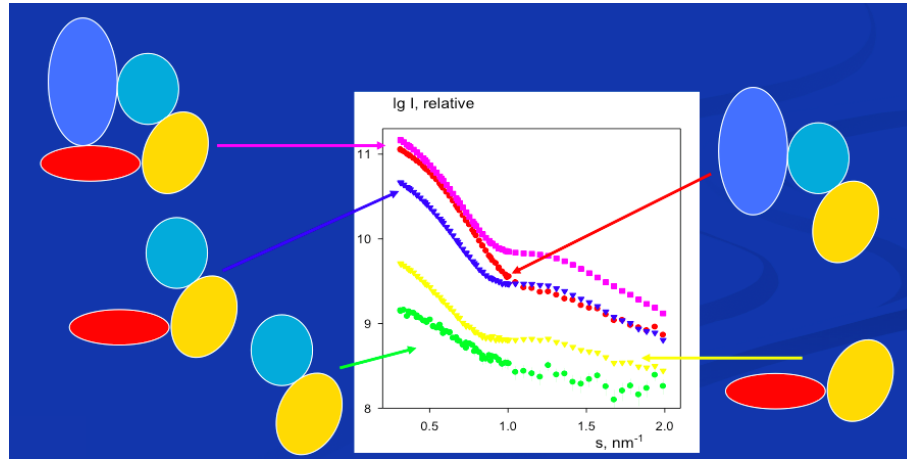
# Rigid and flexible modeling

# Rigid body fitting

- The structures of two subunits in reference positions are known.

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

- This operation depends on three Euler rotation angles and three Cartesian shifts.
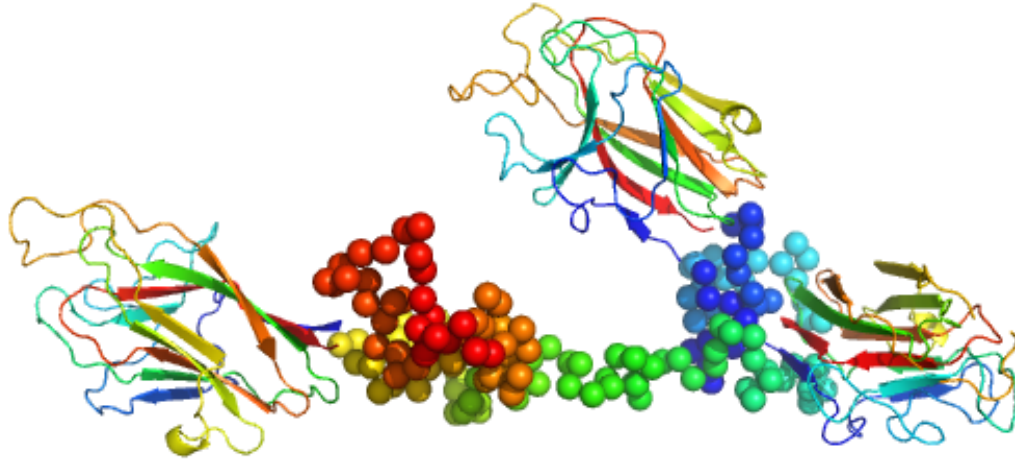
# Rigid body modeling with SASREF



- Fits (multiple X-ray and neutron) scattering curve(s) from partial constructs or contrast variation using simulated annealing

- Requires models of subunits, builds interconnected models without steric clashes.

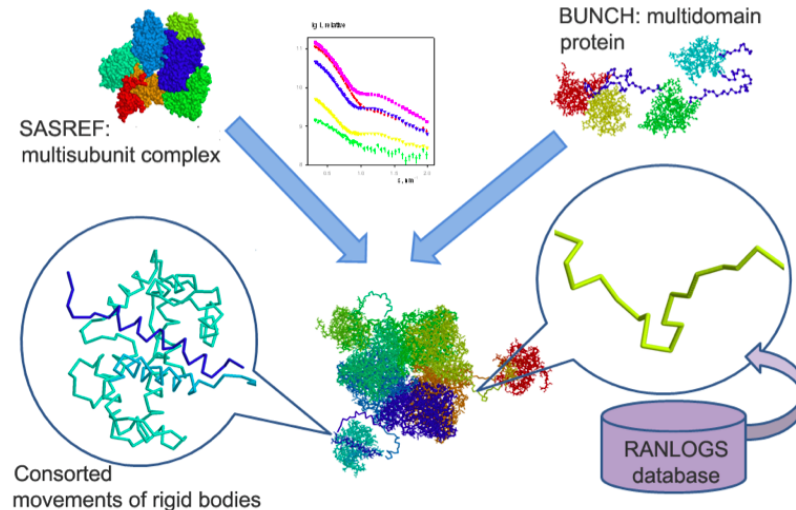- Uses constrains: symmetry, distance, relative orientation if applicable.

Petoukhov & Svergun (2005). *Biophys J.* **89**, 1237; Petoukhov & Svergun (2006). *Biophys J.* **35**, 567

Adapted from Alejandro Panjkovich (EMBL-Hamburg), EMBO 2017

# Addition of missing fragments with BUNCH



- BUNCH combines rigid body and *ab-initio* modelling to find the positions and orientations of rigid domains and probable conformations of flexible linkers represented as dummy residues chains
- Multiple experimental scattering data sets from partial constructs (e.g. deletion mutants) can be fitted simultaneously with the data of the full-length protein.
- accounts for symmetry, allows one to fix some domains and to restrain the model by contacts between specific residues

# Addition of missing fragments with CORAL



- A combination of SASREF and BUNCH to account for missing loops in multi-subunit biological macromolecules.
- Loops are modeled based on known high-resolution structures.

# Summary

- Atsas is a powerful toolbox to analyze SAS data from biological macromolecules

- Dummy models gives an idea of an overall shape of molecule

- Hi-res structures can be compared with SAS data either or used as a building blocks in rigid or flexible modeling

- Potential flaw: you will always get answer but not always correct.

# Questions?