



LUND  
UNIVERSITY



EUROPEAN  
SPALLATION  
SOURCE

# Bayesian inference of conformational ensembles from limited structural data

Wojtek Potrzebowski

Data Management and Software Centre, ESS  
Biochemistry and Structural Biology, LU

[www.europeanspallationsource.se](http://www.europeanspallationsource.se)

7 April, 2019

# European Spallation Source

60% complete

13 funding members, Sweden and Denmark serving as host countries

Pulsed neutron source

More than 100 times brighter neutrons than currently available

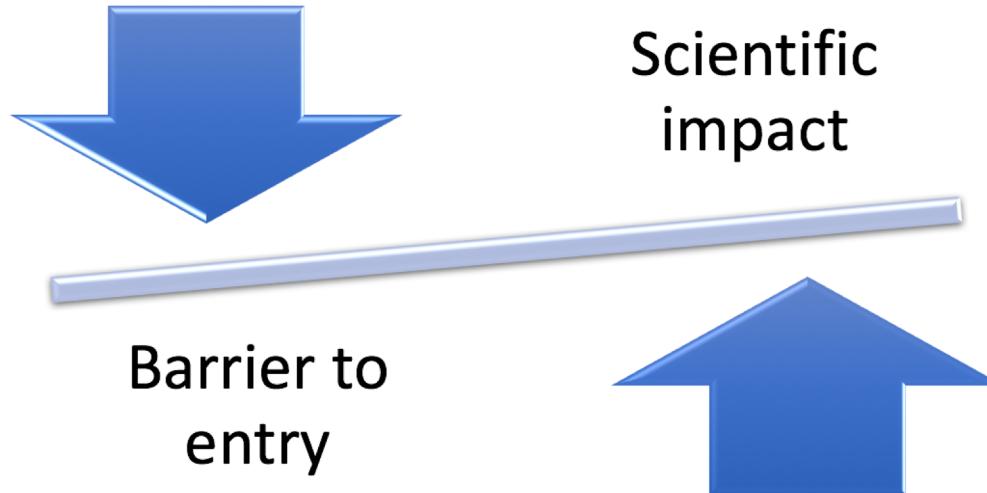
Exponentially more data to collect and analyse



# Data Management and Software Centre (DMSC) balances scientific impact with entry barrier



Minimise the time it takes to understand experimental data



Maximise the scientific impact and success of ESS

**We follow the best software development practice!**

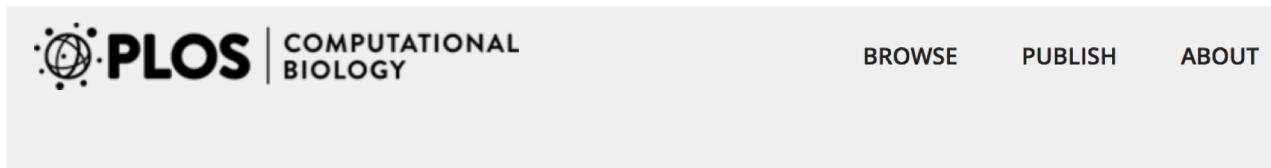
# ESS/DMSC engages in the collaborations



Ingemar André, Lund University



Jill Trehewella, Univ. of Sydney



PLOS COMPUTATIONAL BIOLOGY

BROWSE PUBLISH ABOUT

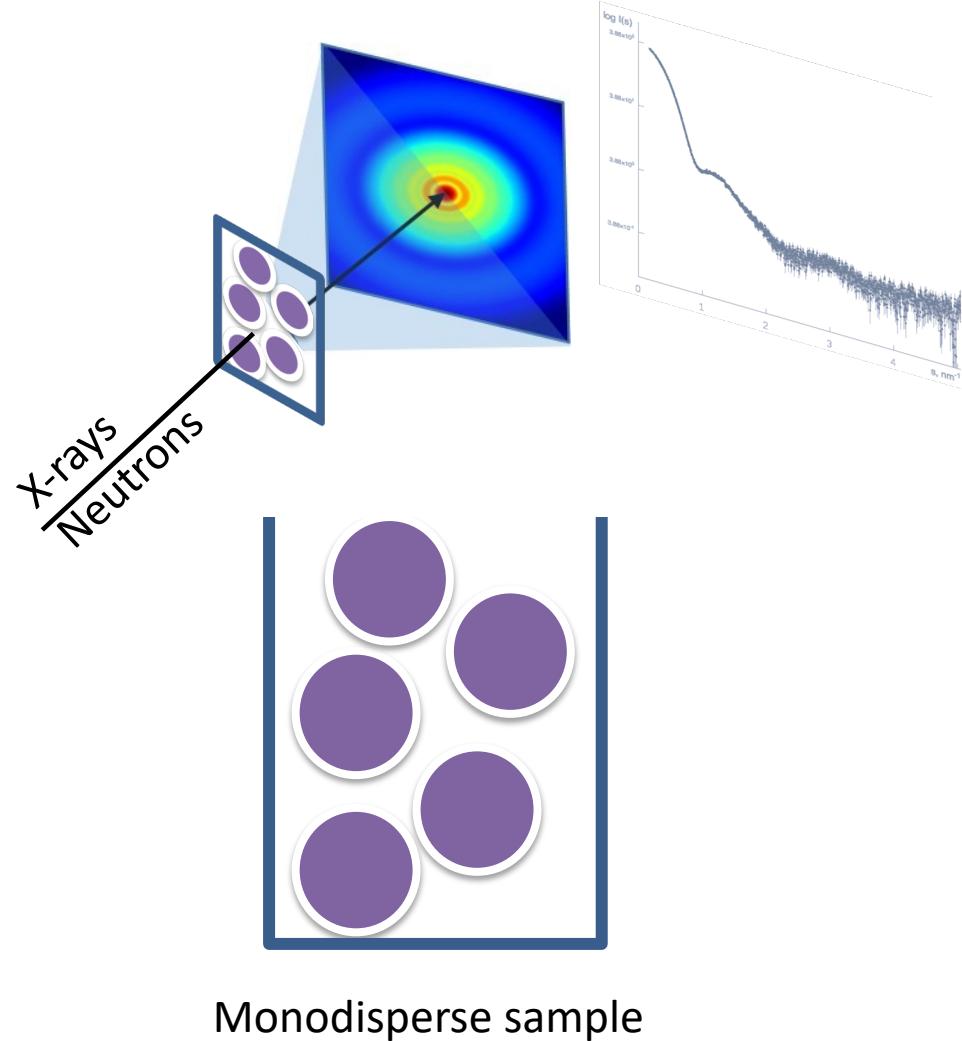
OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

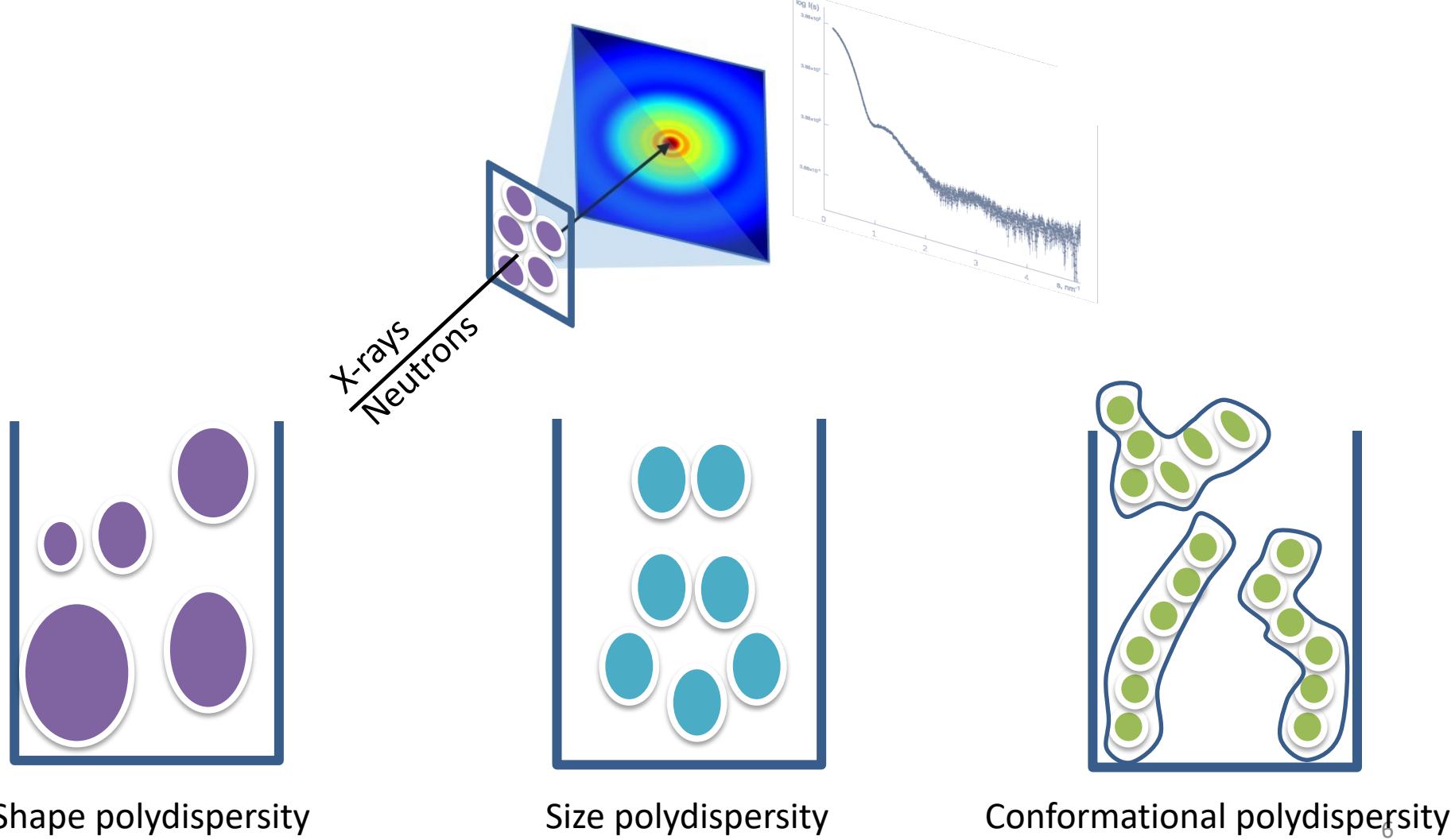
## Bayesian inference of protein conformational ensembles from limited structural data

Wojciech Potrzebowski, Jill Trehewella, Ingemar Andre 

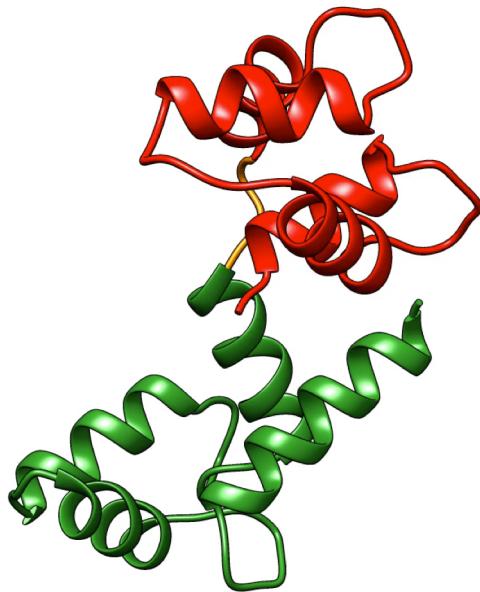
# Monodispersity is desired in SAS experiment



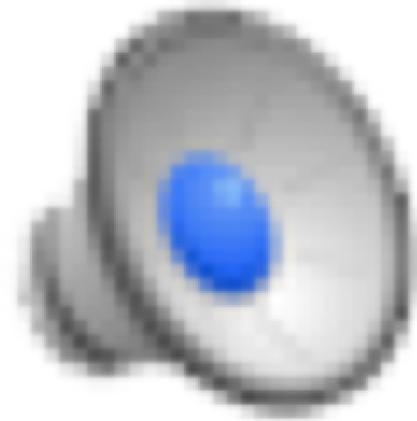
# Polydispersity complicates structural determination



# Proteins with flexible linkers are attractive targets for SAS experiments

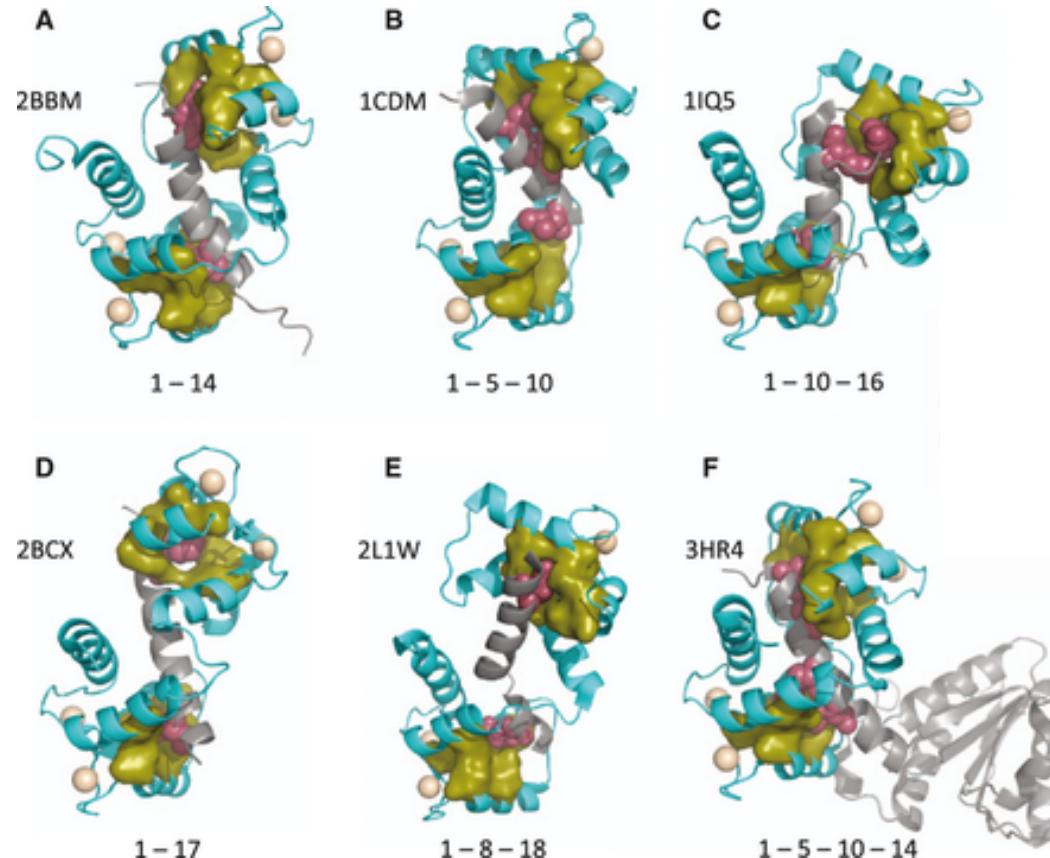


Calmodulin (CaM)



Cardiac myosin binding protein C ( $\Delta mC2$ )

# Domain-domain flexibility plays a central role in biology

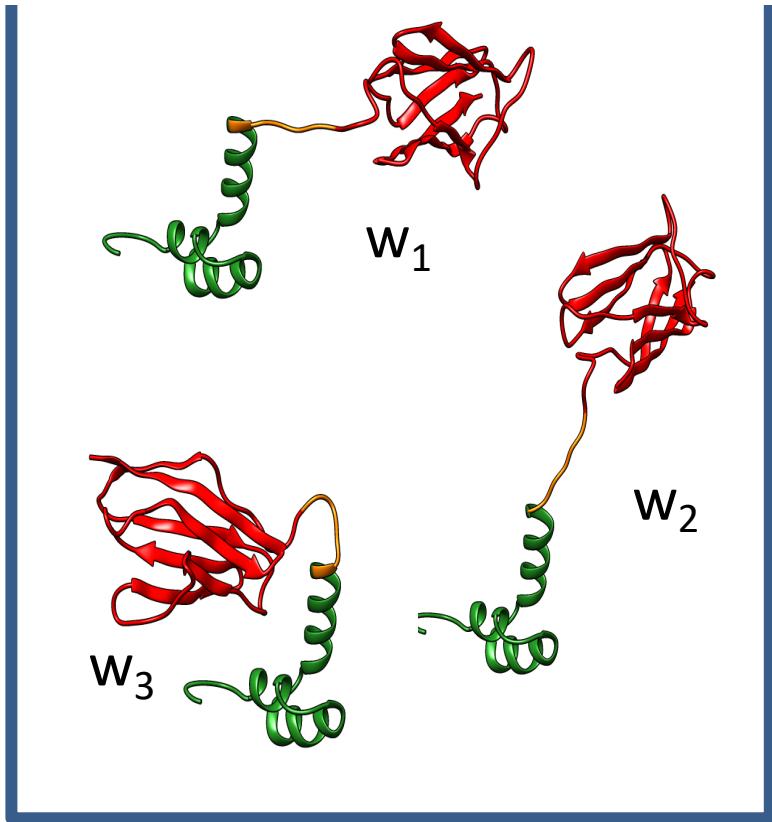


Structural diversity of calmodulin binding to its target sites, Volume: 280, Issue: 21, Pages: 5551-5565, First published: 20 April 2013, DOI: (10.1111/febs.12296)

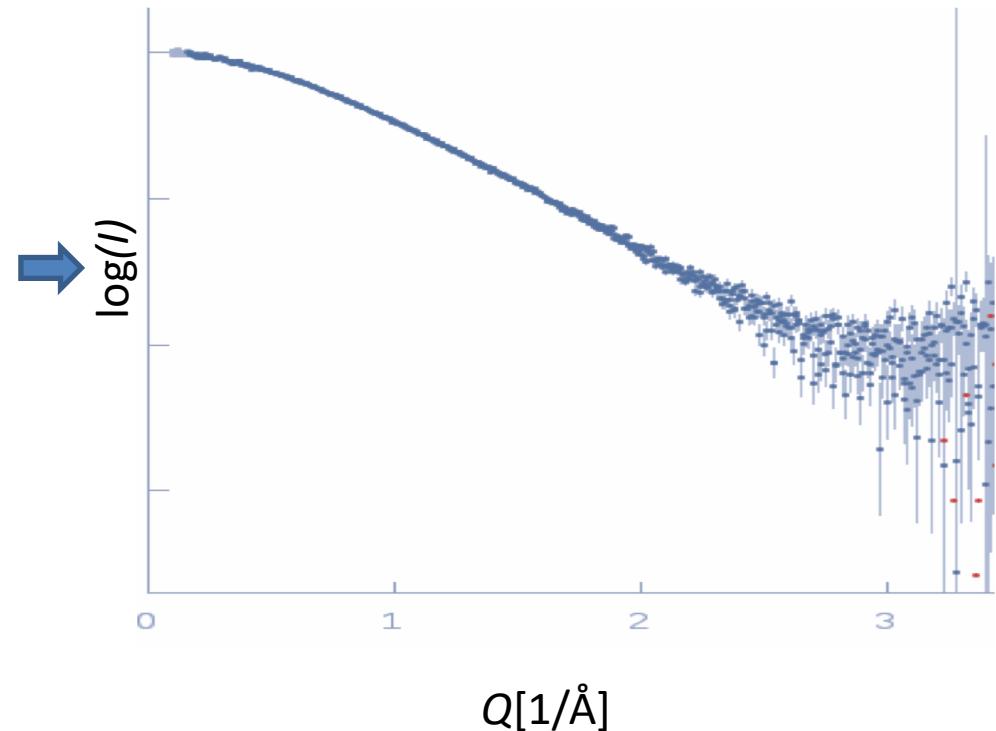
# Challenges for analyzing polydisperse data

- Scattering from the mixture of components
- Low information content
- Overfitting to experimental noise
- Optimization problem – thousands of parameters

# Multiple conformations contribute to scattering curve

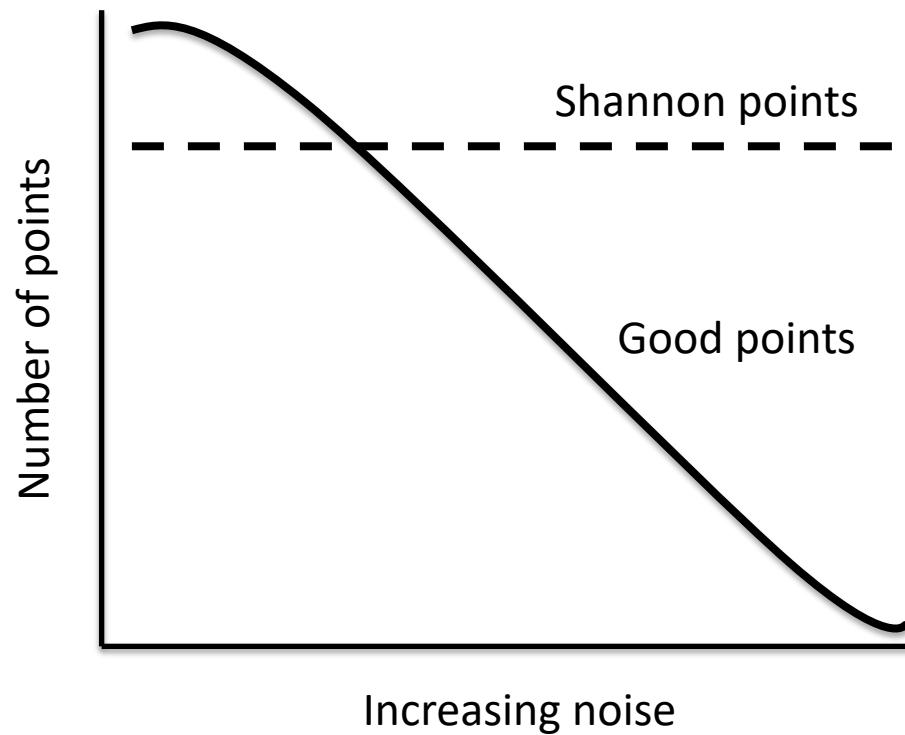


$$w_1 + w_2 + w_3 = 1$$



# SAXS data contains limited structural information

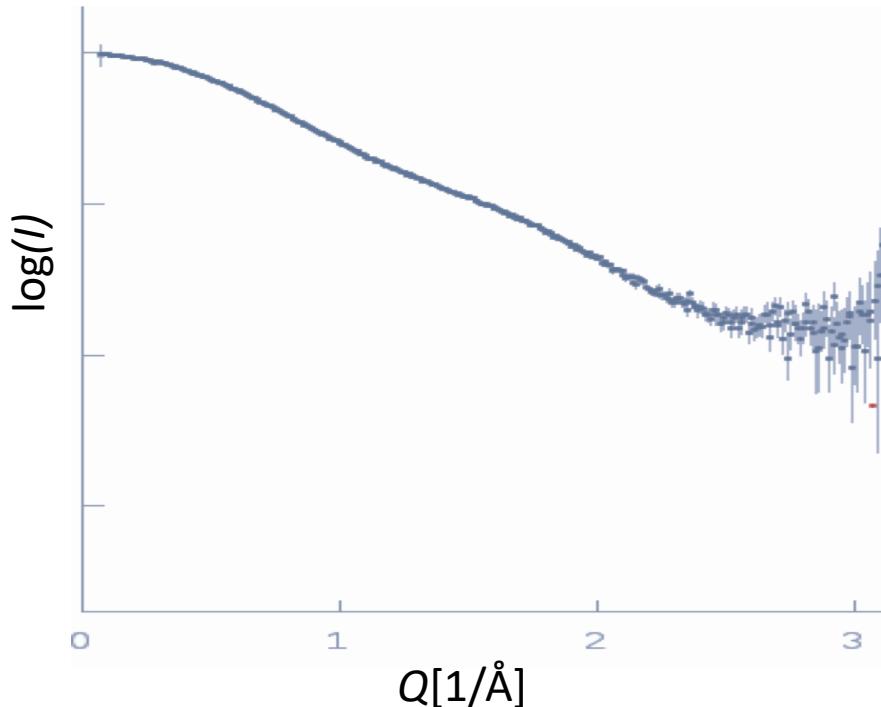
The number of independent data points estimated from Nyquist-Shannon sampling theorem or maximum entropy regularization (Vestergaard and Hansen (2006))



The number of independent measurement rarely exceeds **5-15**.

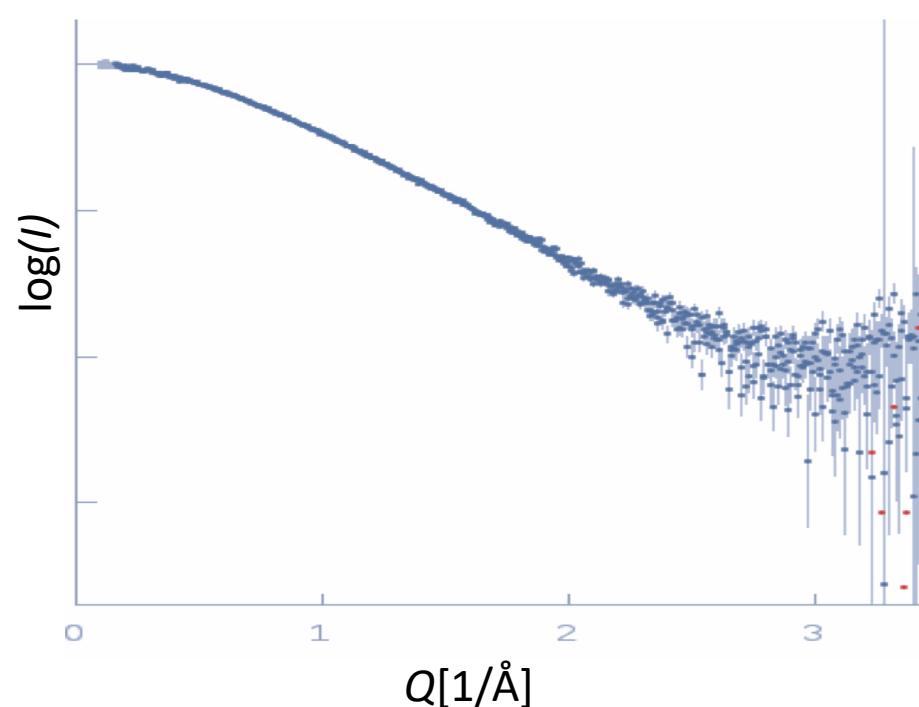
# SAXS data contains limited structural information

Calmodulin (SEC-SAXS)



$$N_s = 6.4$$
$$N_g = 5.1$$

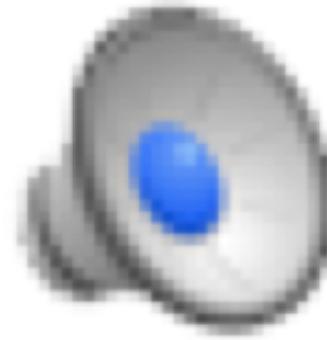
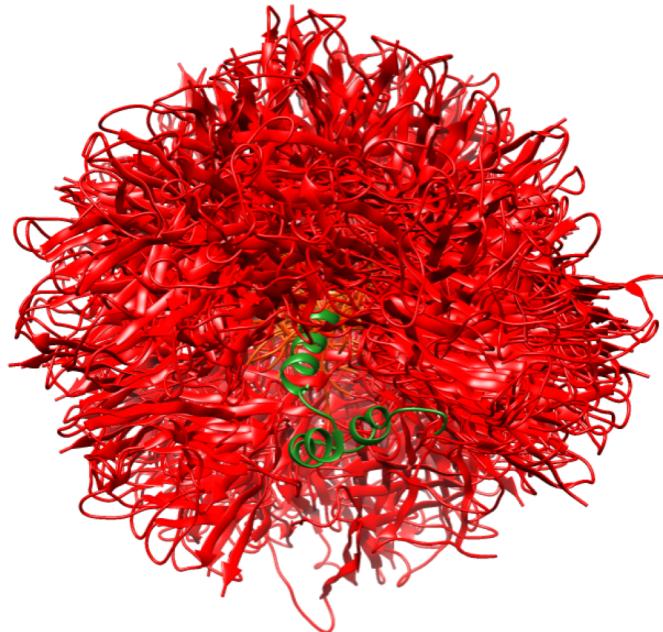
Cardiac myosin-binding protein C (SAXS)



$$N_s = 8.7$$
$$N_g = 6.1$$

# Many parameters describe conformational ensembles

- Thousands of parameters may be required to define ensemble
- Ensembles can be represented in discrete or continuous space



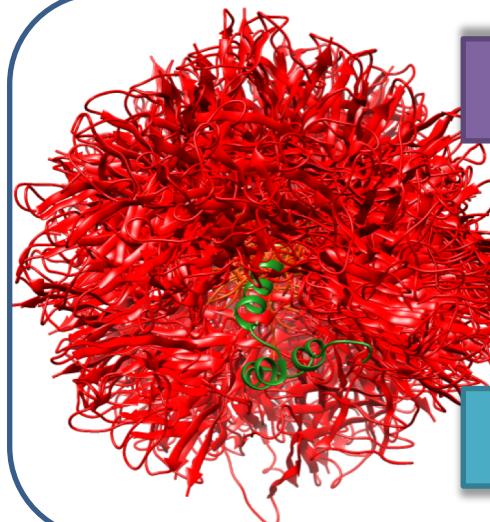
- $\chi^2$  based optimization leads to overfitting

# What do we want from our inference method?

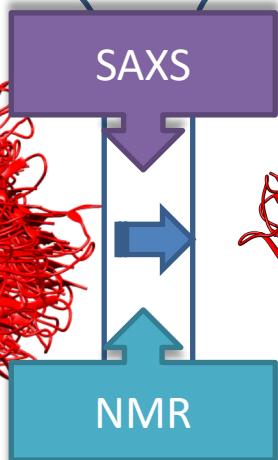
- We want to extract as much information from the experimental data – but not more than that
  - Avoid overfitting to experimental noise
  - Quantify the uncertainty in the inference
  - Combine multiple sources of data: simulation and experimental data sources
- > Bayesian Inference to the rescue!

# Method to infer conformational ensembles

Monte Carlo simulations



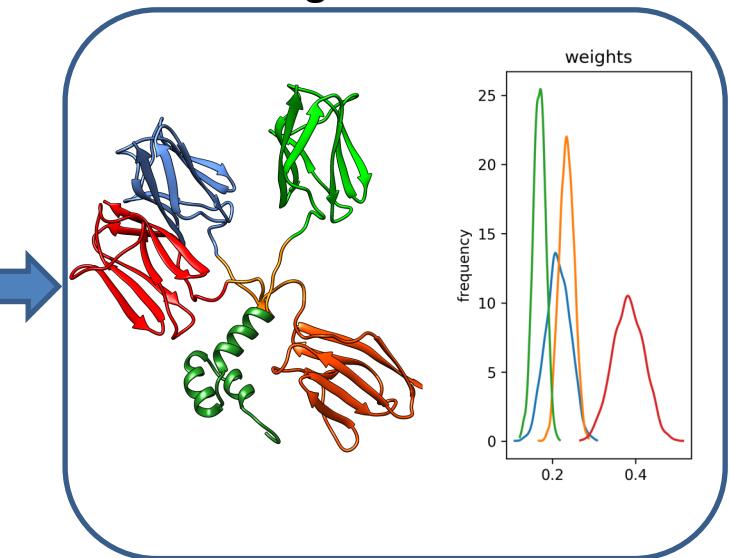
Model selection



Approximate Bayesian



Weights inference

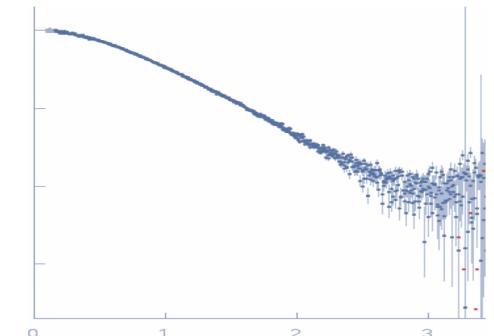
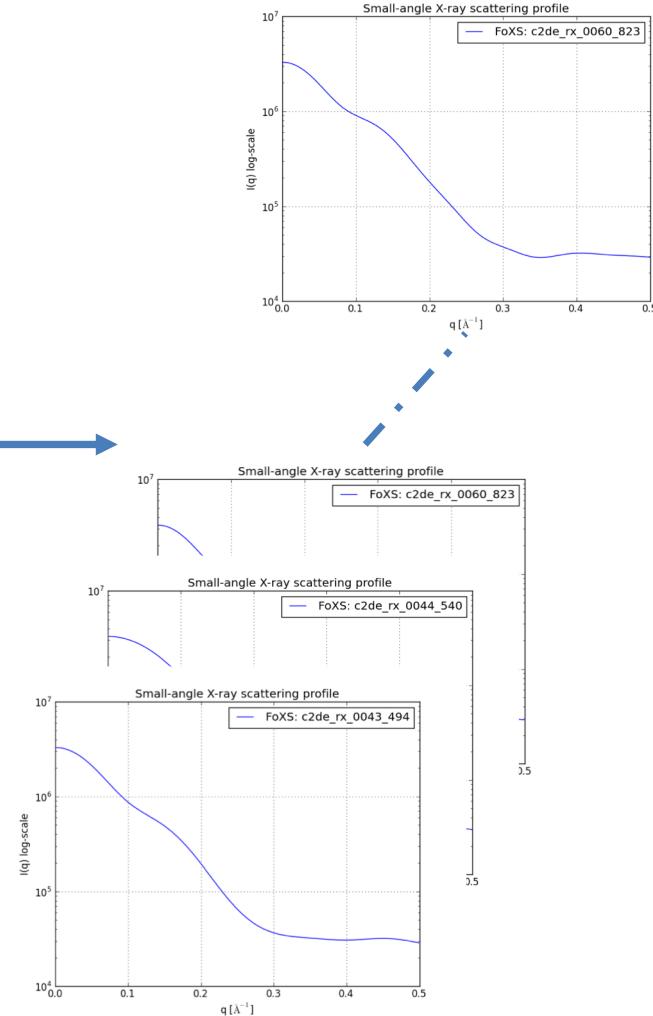
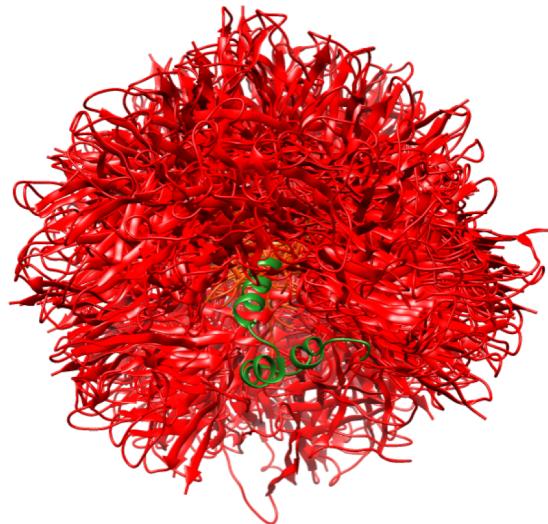


Rosetta structure prediction



Complete Bayesian

# Scattering patterns are generated after Rosetta modeling (once)



Experimental

# Bayesian statistics can be applied to infer weights

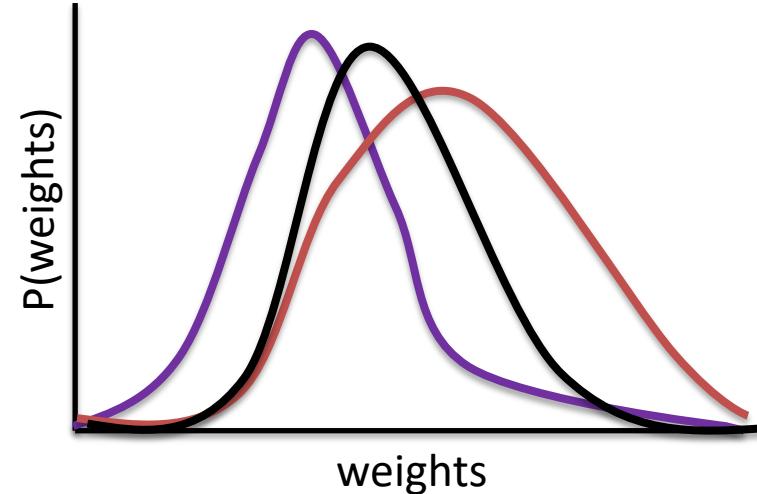
**Bayes rule:**



Thomas Bayes

$$P(w|Data) = \frac{P(Data|w) * P(w)}{P(Data)}$$

*Posterior  $\propto$  likelihood \* prior*



Posterior ( $P(w|Data)$ ): Probability of  $w$  after observing data

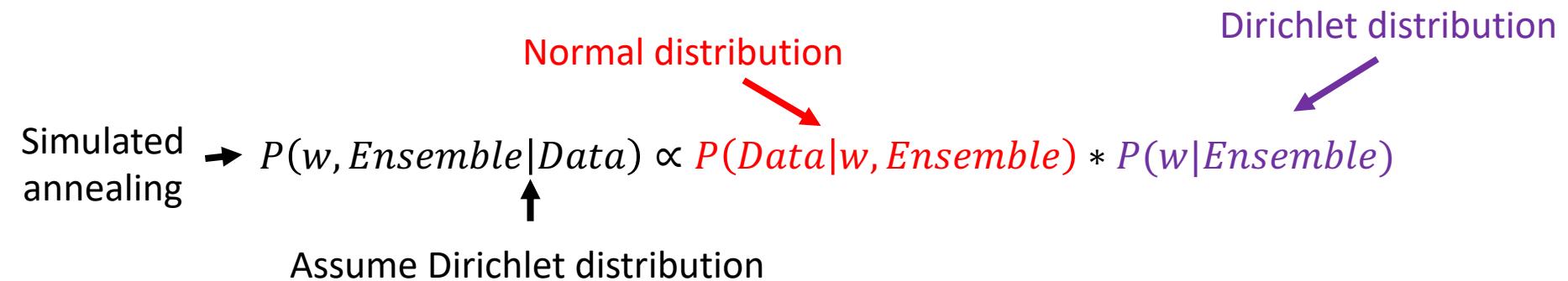
Prior ( $P(w)$ ): Information before observation of data

Likelihood ( $P(Data|w)$ ): Probability of observing data given  $w$

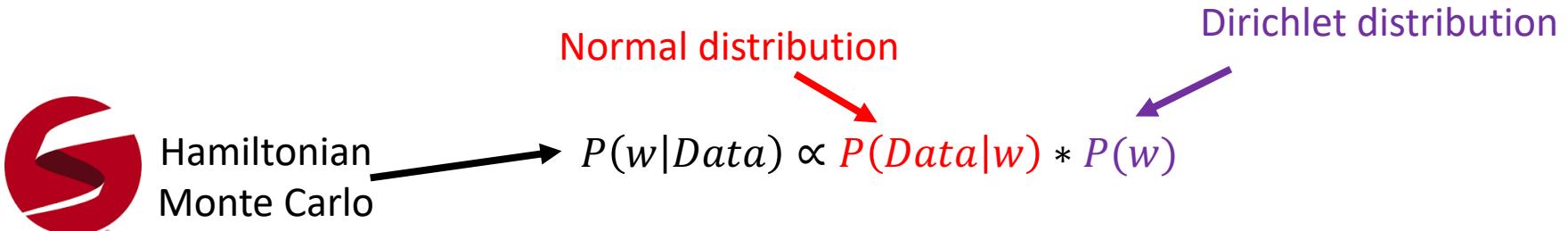
$P(Data)$ : Normalization constant also referred to as Model evidence

# Approximate and complete inference complements each other

Approximate Bayesian inference – fast, less accurate:

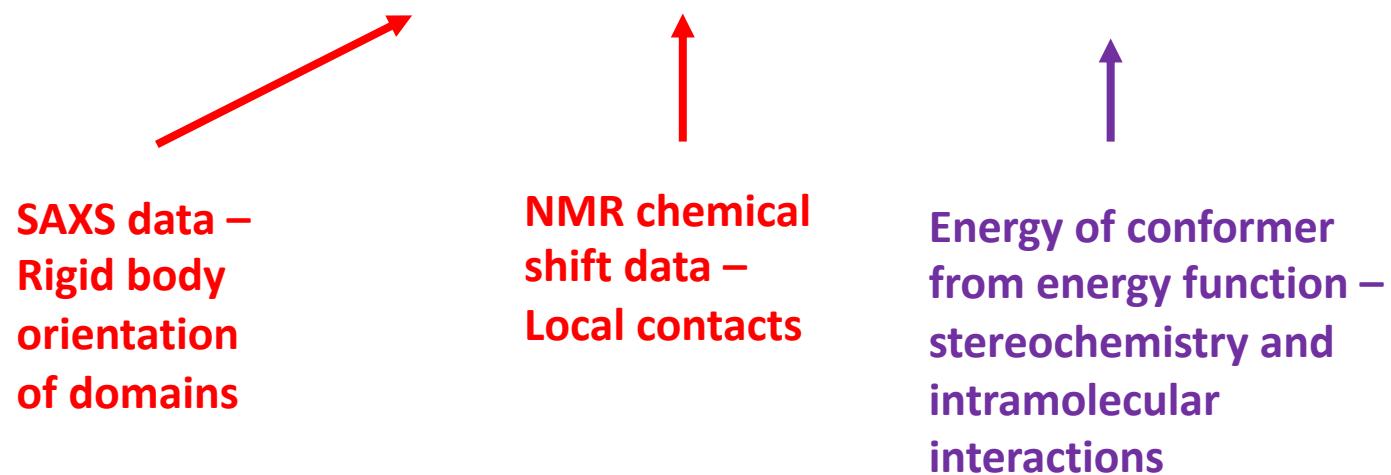


Complete Bayesian inference - slow:



# Combining multiple data sources increases structural information

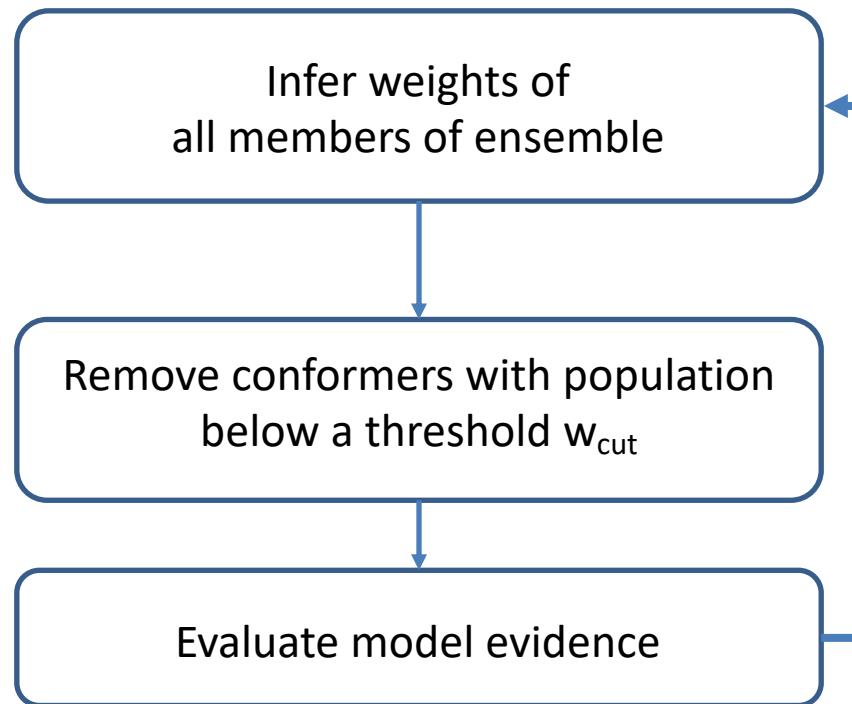
$$P(w|SAXS, NMR, \text{conformer energy}) \propto P(SAXS|w) * P(NMR|w) * P(w(\text{conformer energy}))$$



# Approximate Bayesian inference balances complexity with fit to the data

For a given set of conformers

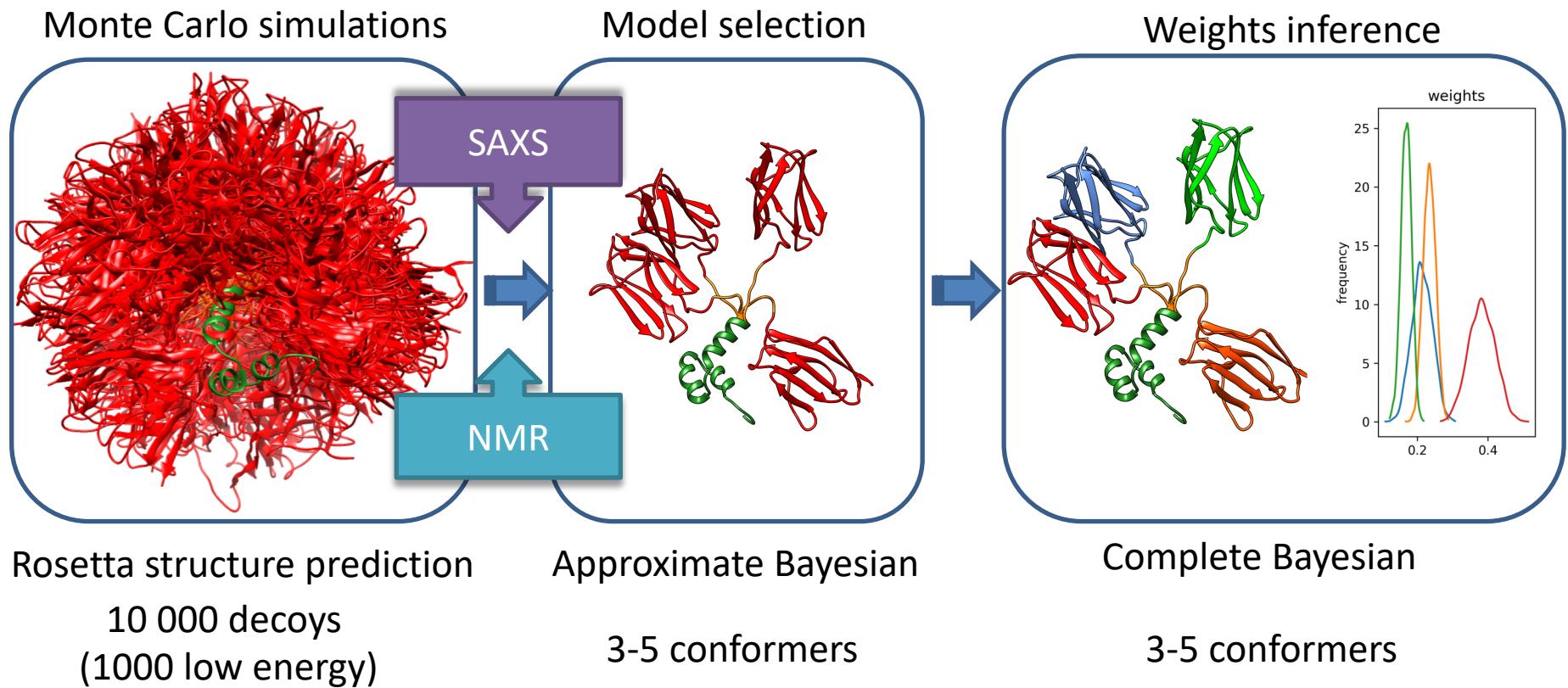
Fisher CK, Ullman O, Stultz C. (2012 )  
Pacific symposia on biocomputing.



$$P(w, \text{Ensemble} | \text{Data}) = \frac{P(\text{Data} | w, \text{Ensemble}) * P(w | \text{Ensemble})}{P(\text{Data} | \text{Ensemble})}$$

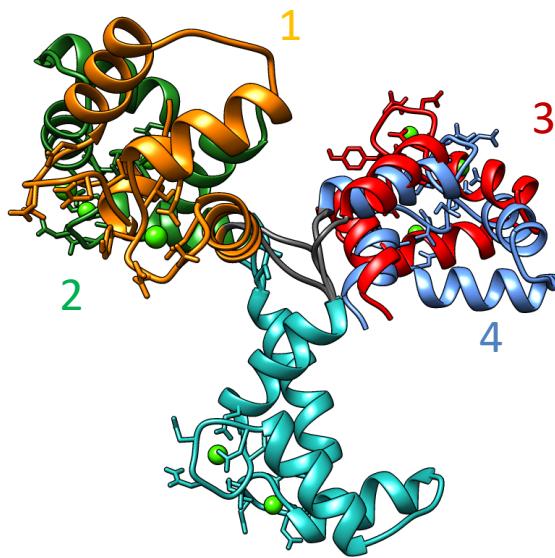
**Model evidence** —→  $P(\text{Data} | \text{Ensemble}) \propto P(\text{Ensemble} | \text{Data})$

# Method to infer conformational ensembles

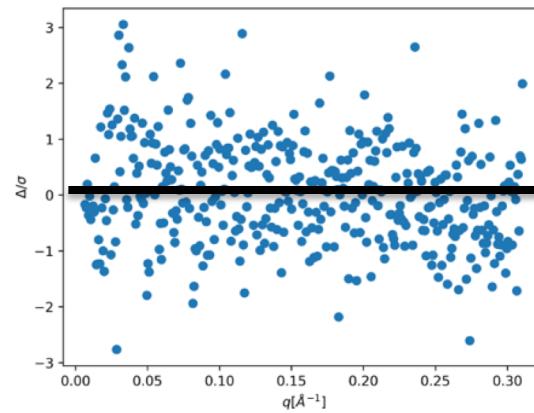
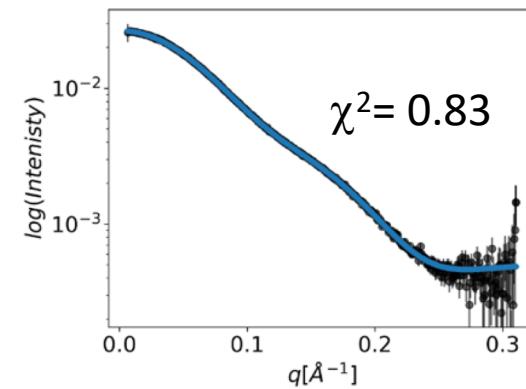
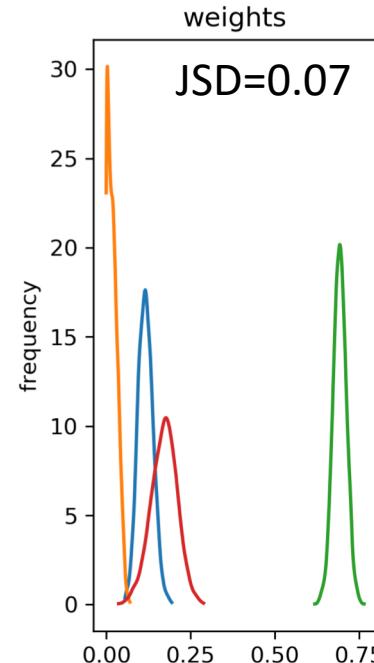


# Calmodulin ensemble inferred from SAXS data

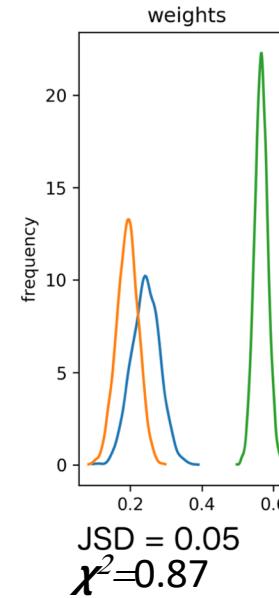
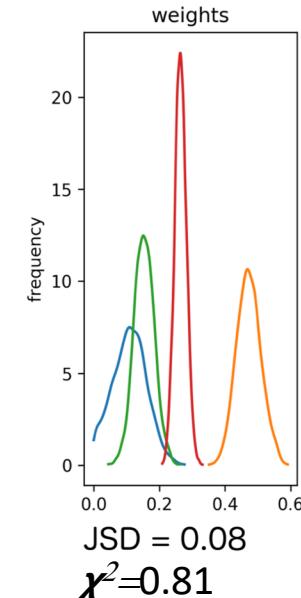
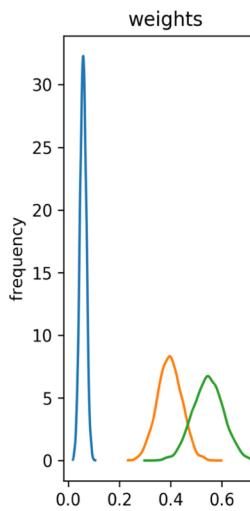
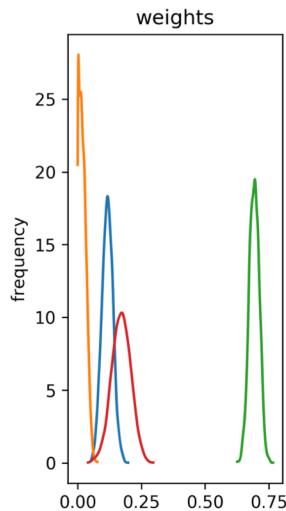
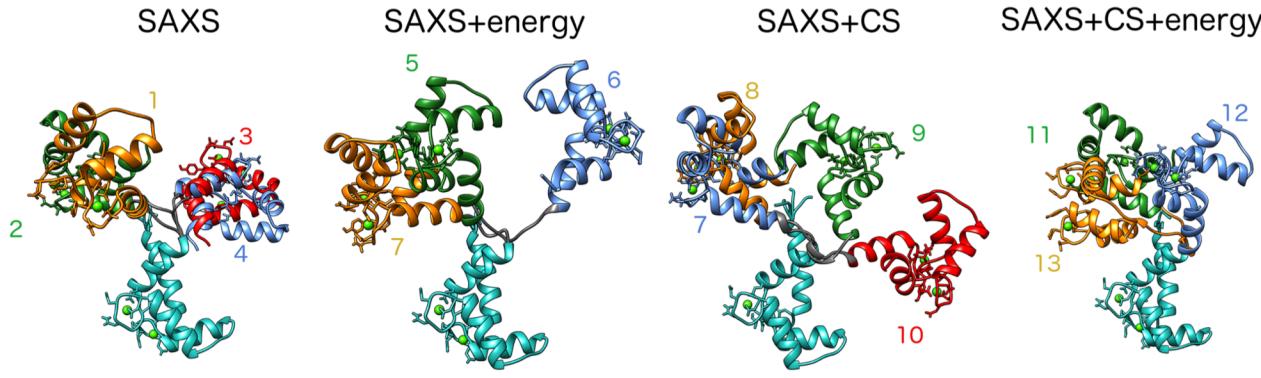
$$\text{posterior} \propto \text{likelihood(SAXS)} \cdot \text{prior(weights)}$$



Conformer	Population weight
1	$0.02 \pm 0.01$
2	$0.69 \pm 0.02$
3	$0.17 \pm 0.04$
4	$0.12 \pm 0.02$

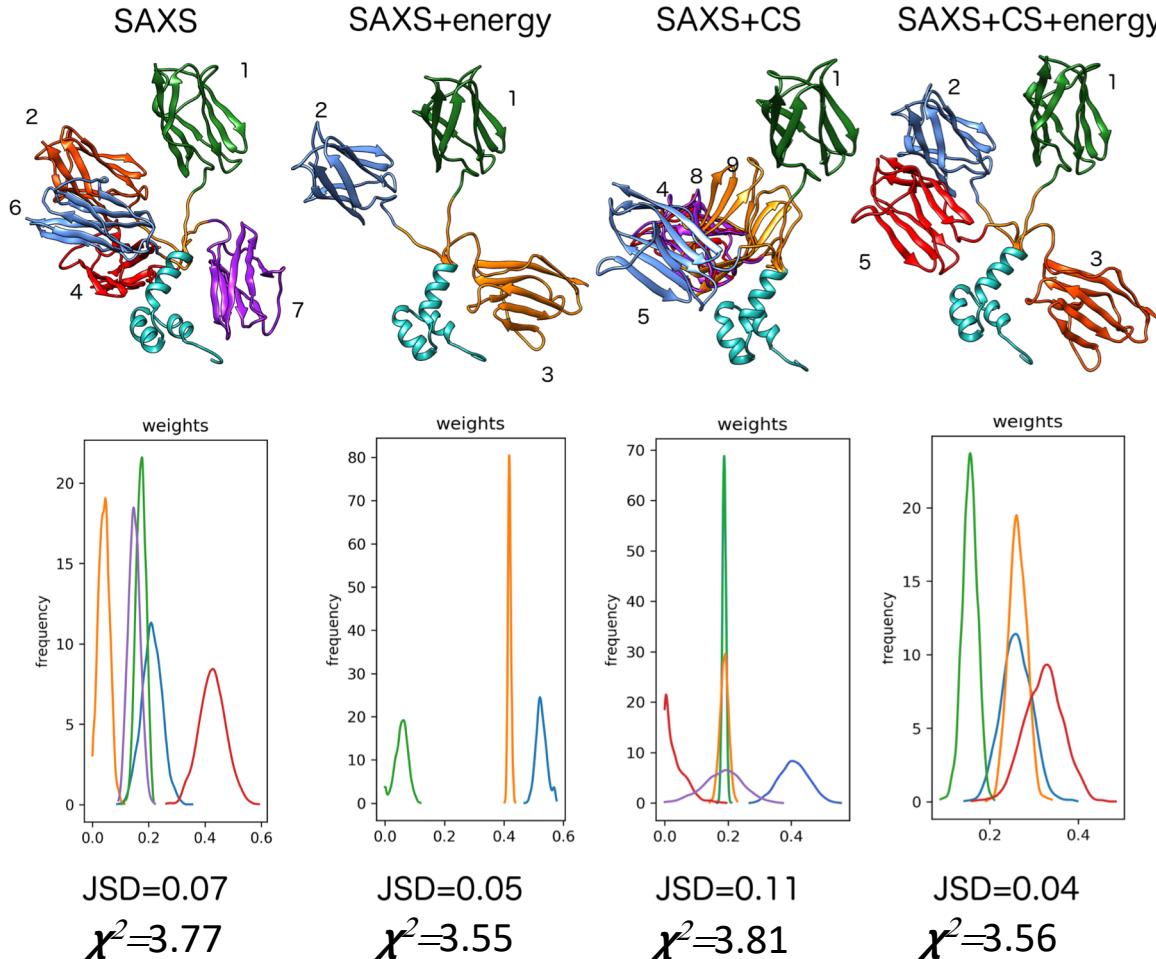


# Combination of different data sources explains SAXS data equally well



energy = energy from molecular simulations  
CS = NMR chemical shifts  
JSD = Jensen-Shannon weights divergence

# Combination of different data sources improves definition of $\Delta mC2$ ensemble



energy = energy from molecular simulations  
CS = NMR chemical shifts  
JSD = Jensen-Shannon weights divergence

*The drop in  $\chi^2$  is the result in the improved quality of ensemble*

# Conclusions I

## Bayesian Inference:

- can be used to optimally balance model complexity and fit to experimental data
- gives information about the uncertainty of inferred ensembles and parameters
- allows for combining information from different data sources in one model

Our method combines the speed of alternative methods with the benefits of Bayesian inference

# Method is available for download

<https://andre-lab.github.io/bioce/>

**bioce**  
bayesian inference of conformational ensembles

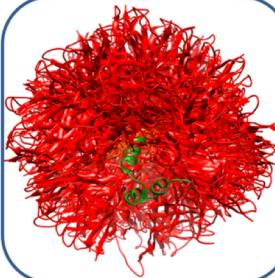
[View On GitHub](#) 



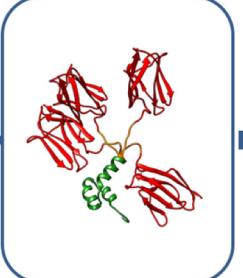
## Method overview

A method based on Bayesian statistics that infers conformational ensembles from a structural library generated by all-atom Monte Carlo simulations. The first stage of the method involves a fast model selection approach based on variational Bayesian inference that maximizes the model evidence of the selected ensemble. This is followed by a complete Bayesian inference of population weights in the selected ensemble.

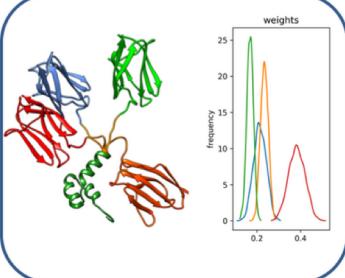
**Rosetta structural models**



**Model selection**



**Weights inference**



weights

Weight Range	Frequency
0.15 - 0.20	~12
0.20 - 0.25	~22
0.25 - 0.30	~15
0.30 - 0.35	~10
0.35 - 0.40	~12

Approximate Bayesian

Full Bayesian

- [Method overview](#)
- [Installation](#)
- [Running examples](#)
- [Generating input data](#)
- [Output](#)
- [Using chemical shift data](#)
- [Using structural energies](#)
- [Webserver](#)

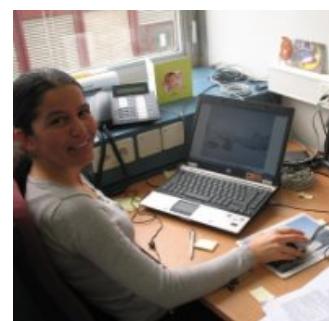
# Exploring protein association pathways with time resolved SANS/SAXS



Ingemar André  
(LU)



Ryan Olivier  
(LU)



Najet Mahmoudi  
(ISIS)



Thomas Holm Rod  
(ESS)



Martin Nors Pedersen  
(KU)



Lise Arleth  
(KU)

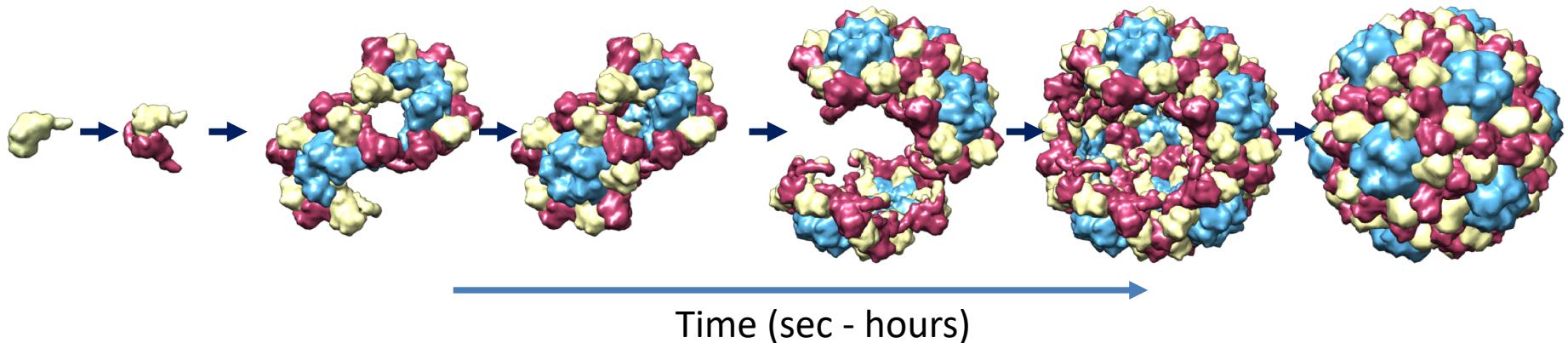


EUROPEAN  
SPALLATION  
SOURCE



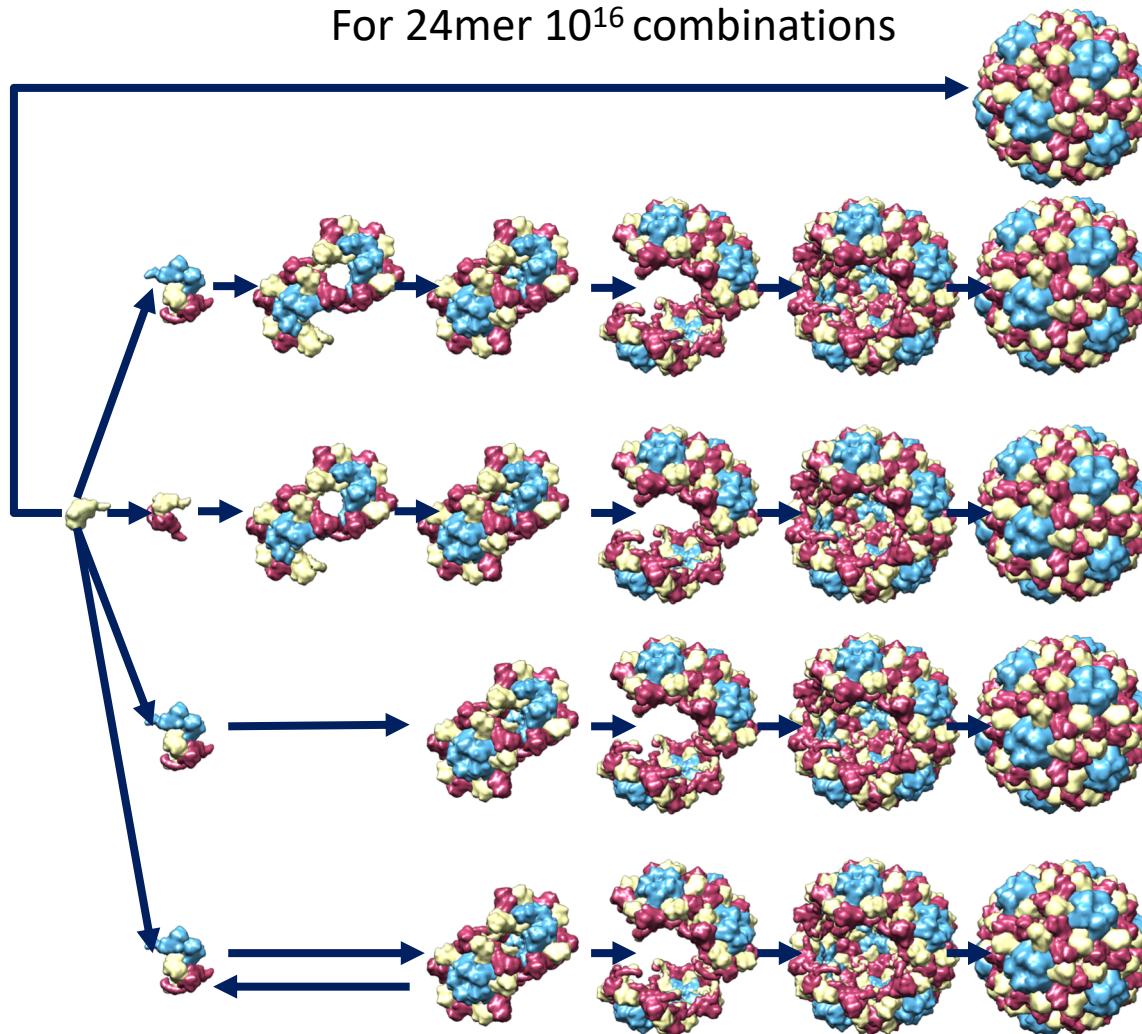
# Why to study protein association pathways?

- Self-assembly central to many processes in different domains
- Key to understanding life processes at molecular level
- For virus capsids:
  - Start and end points are often known
  - But not what happens in between
  - E.g. critical in designing antiviral therapies

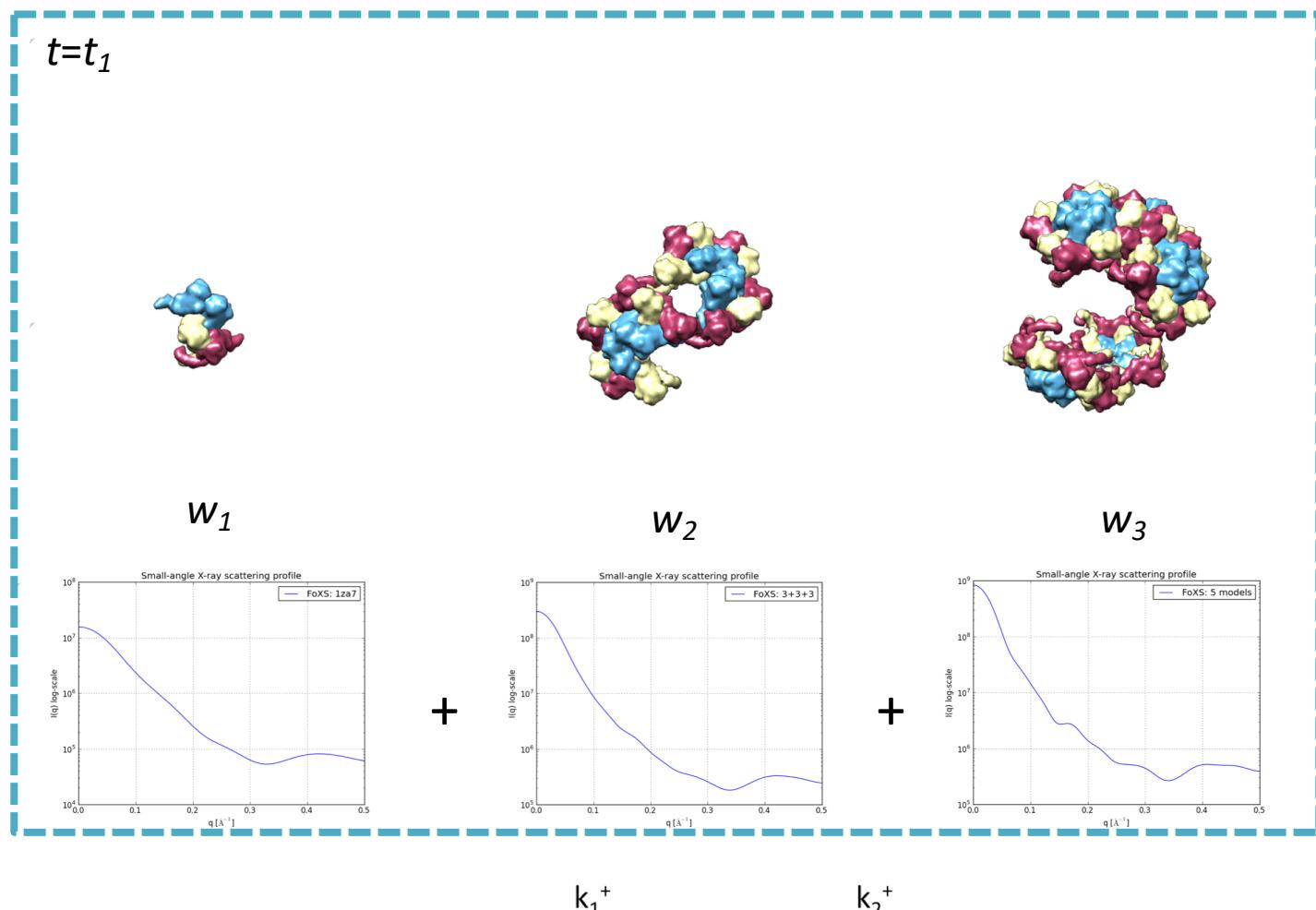
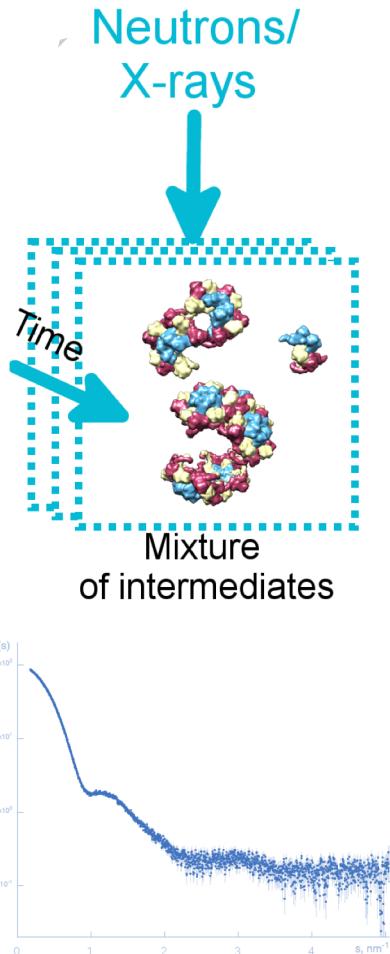


Cowpea chlorotic mottle virus association pathway?

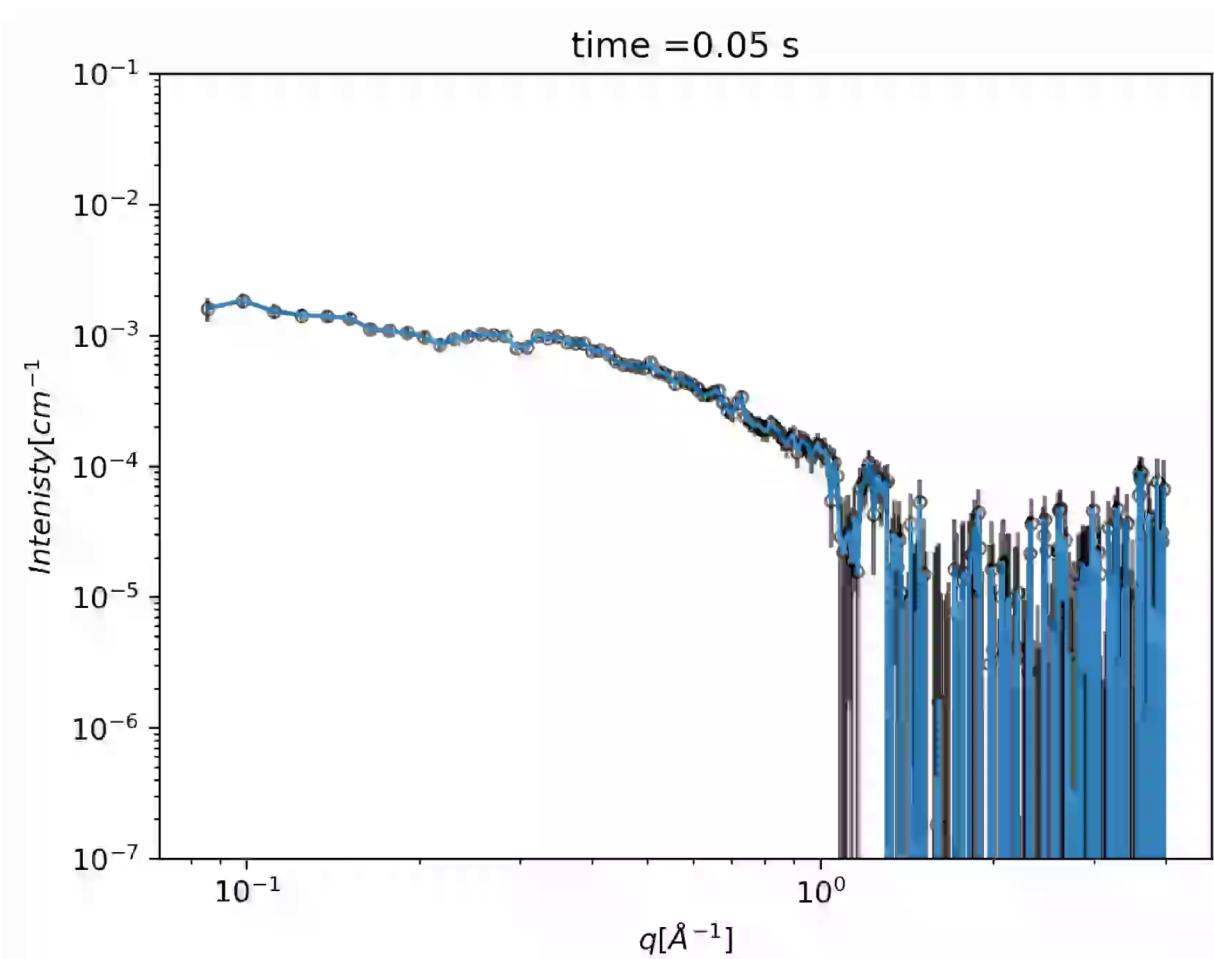
# Humongous number of possible pathways



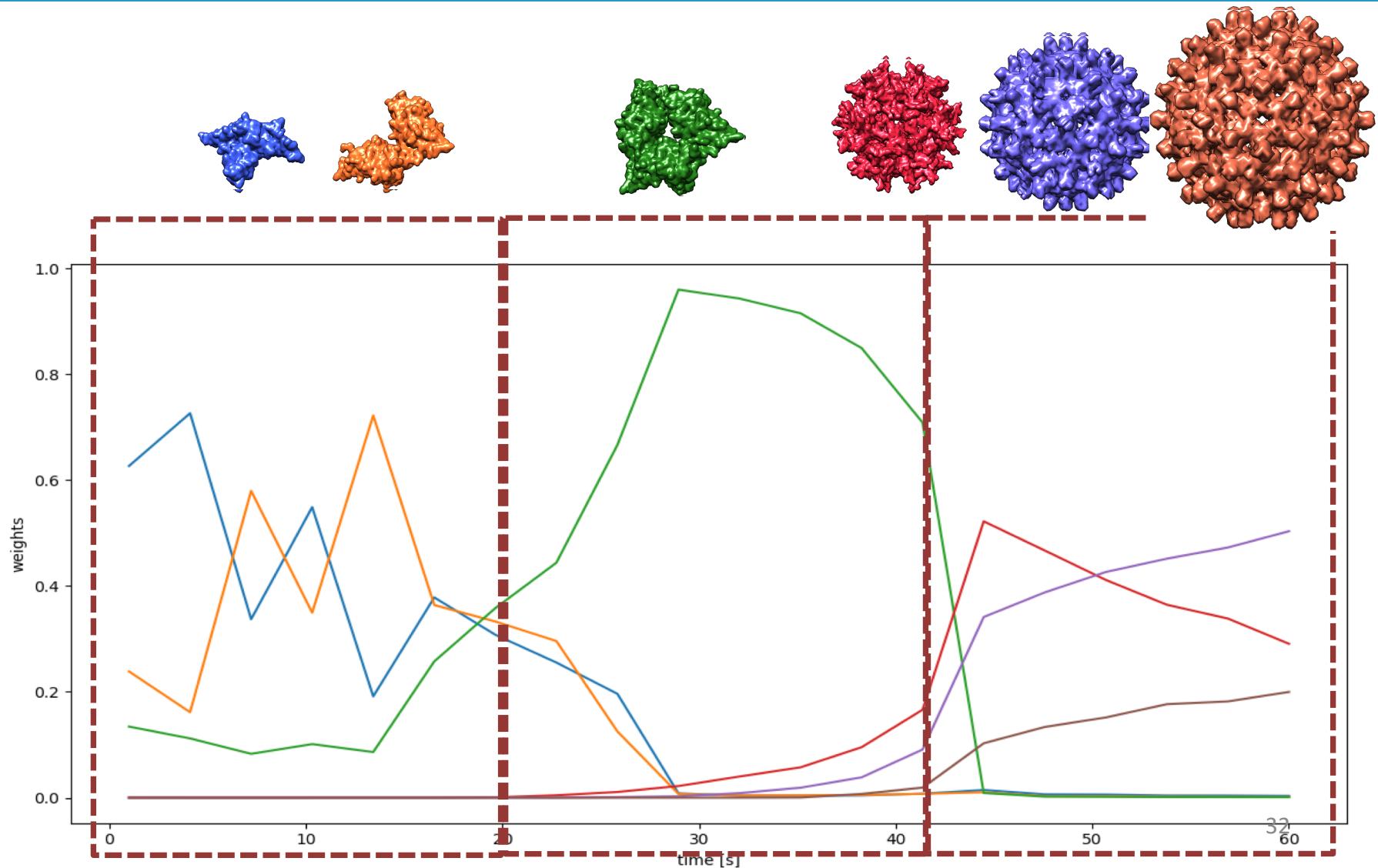
# Method for analyzing SAS data from the mixture of components



# Time-resolved SAXS captures virus capsid assembly

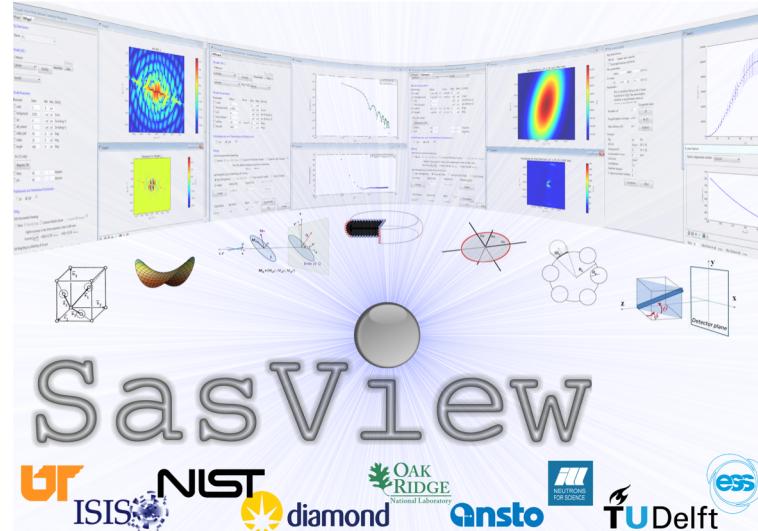


# Intermediates can be inferred from time-resolved SAXS data

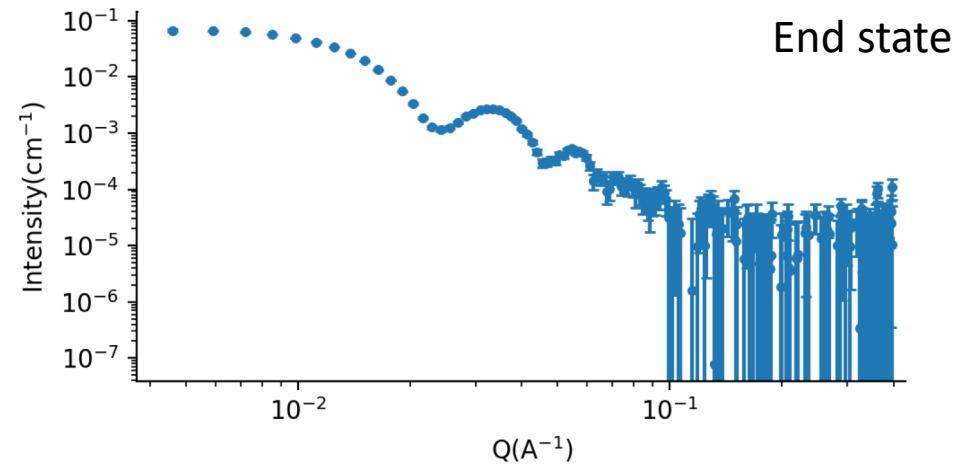
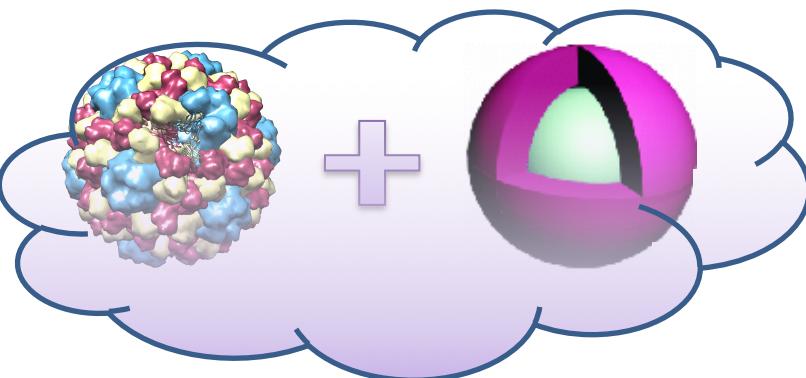


# SasView analyzes and models SAS data

- Community driven project
- 30 developers from X-ray and neutron facilities
- Engaged user community
- Regular code camps and releases
- **We are always on the lookout for new collaborations**
  
- Fitting geometrical models to data
- P(r) inversion
- Model-free calculations
- Correlation functions
- Numerous utility tools



# Bayesian inference combines atomistic and coarse-grained models



Parameter	Value
Weight T3 capsid	0.4
Weight T4 capsid	0.15
Weight core shell sphere	0.45
Core radius [Å]	90
Shell radius [Å]	30

# Conclusions II

- Bayesian inference methods can be adapted to analyze time-resolved SAXS data
- Models of various accuracy can be combined into a single statistical framework

# Thank you!

Lund University:

Ingemar André  
Ryan Oliver

ESS:

Thomas Holm Rod

Univ. of Sydney:

Jill Trewella

ISIS neutron and muon source:

Najet Mahmoudi

Univ. of Copenhagen:

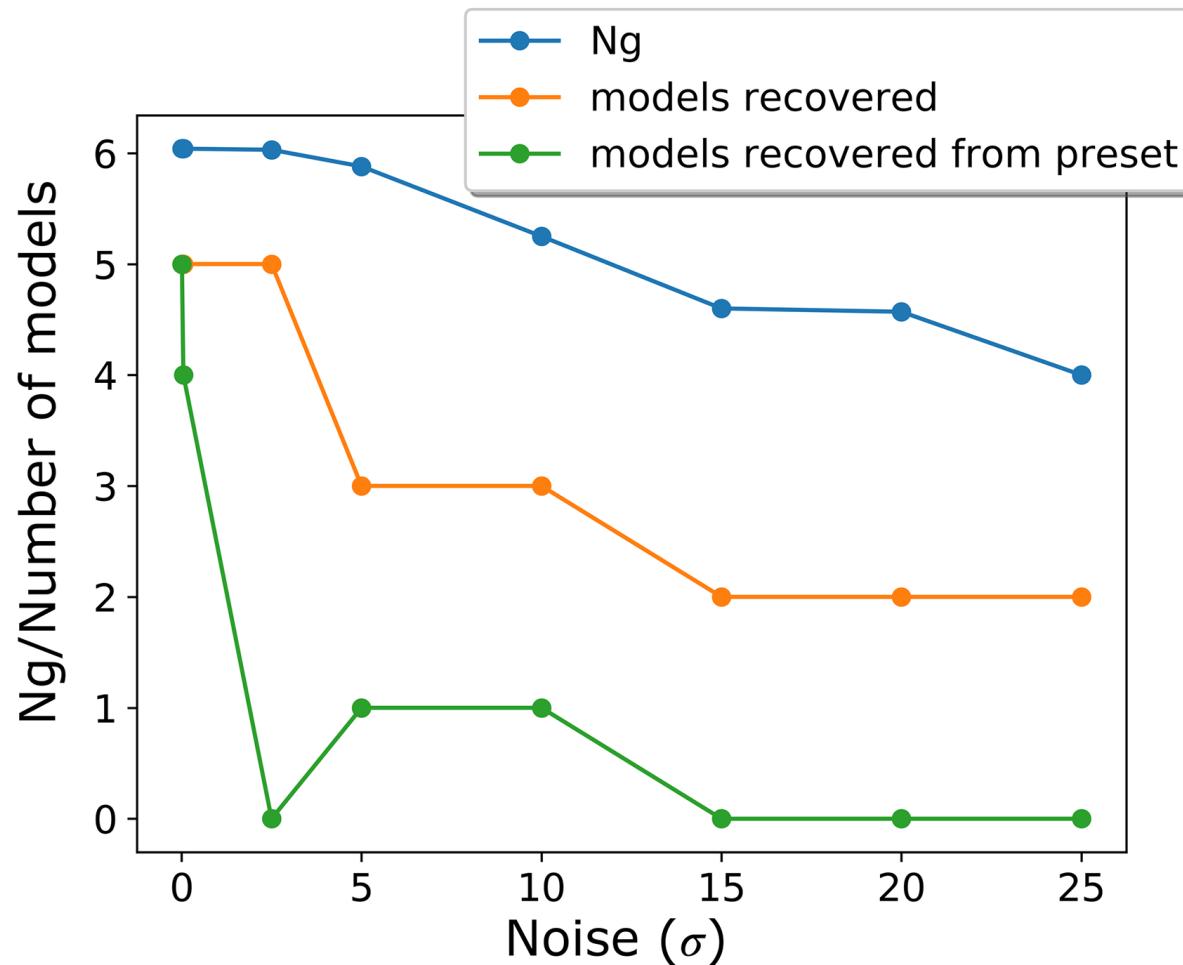
Martin Nors Pedersen  
Lise Arleth

SasView collaboration

Funding:

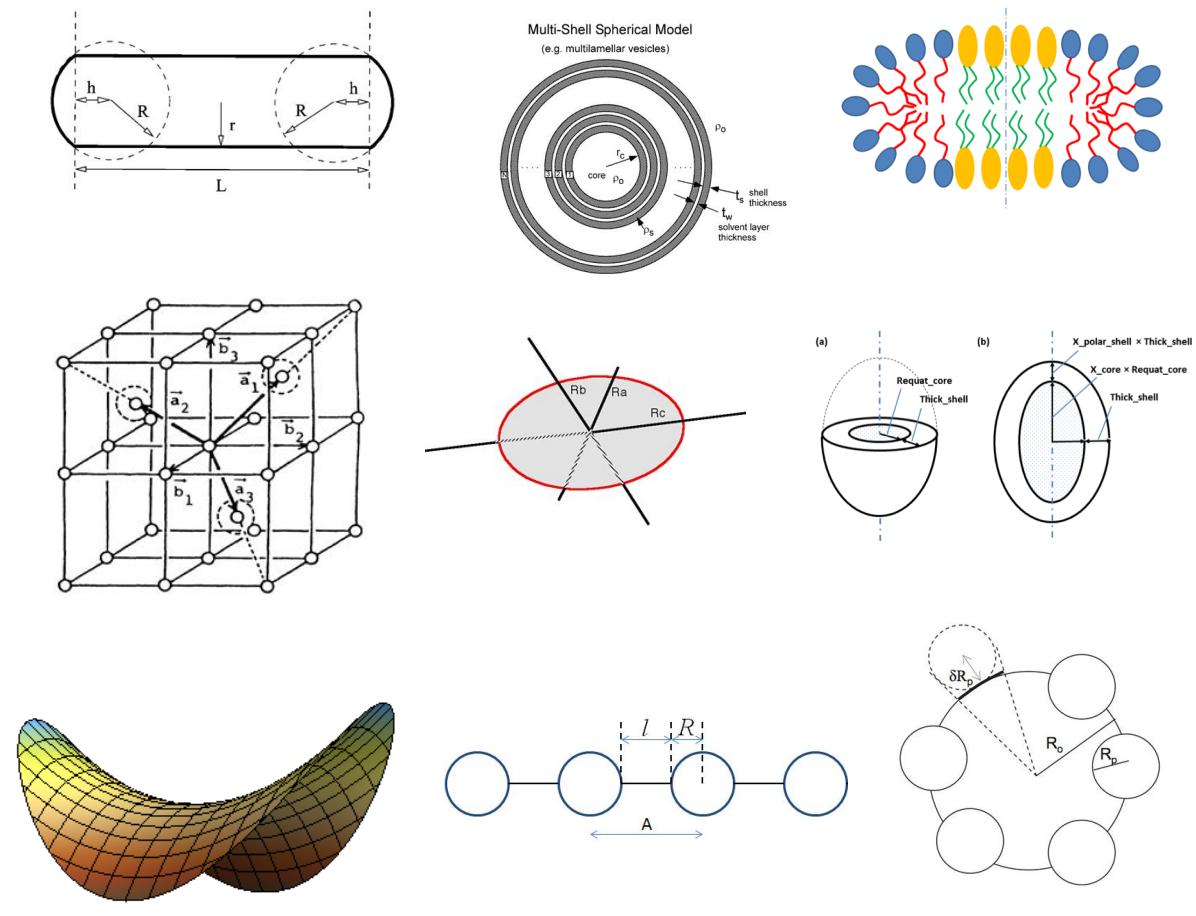
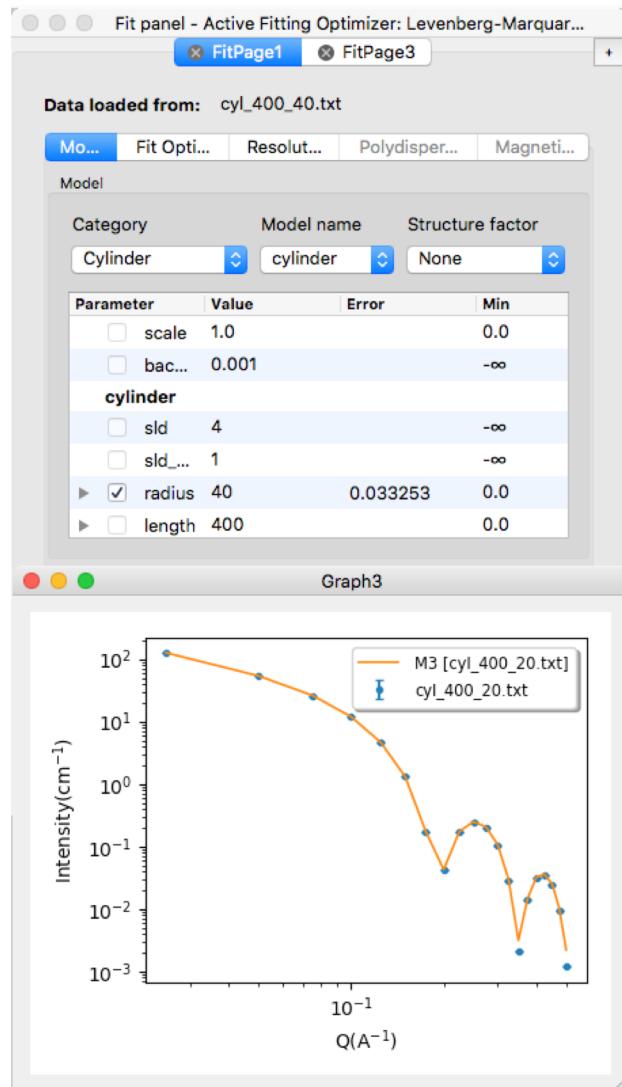


# Approximate bayesian recovers ensembles up to high level of noise

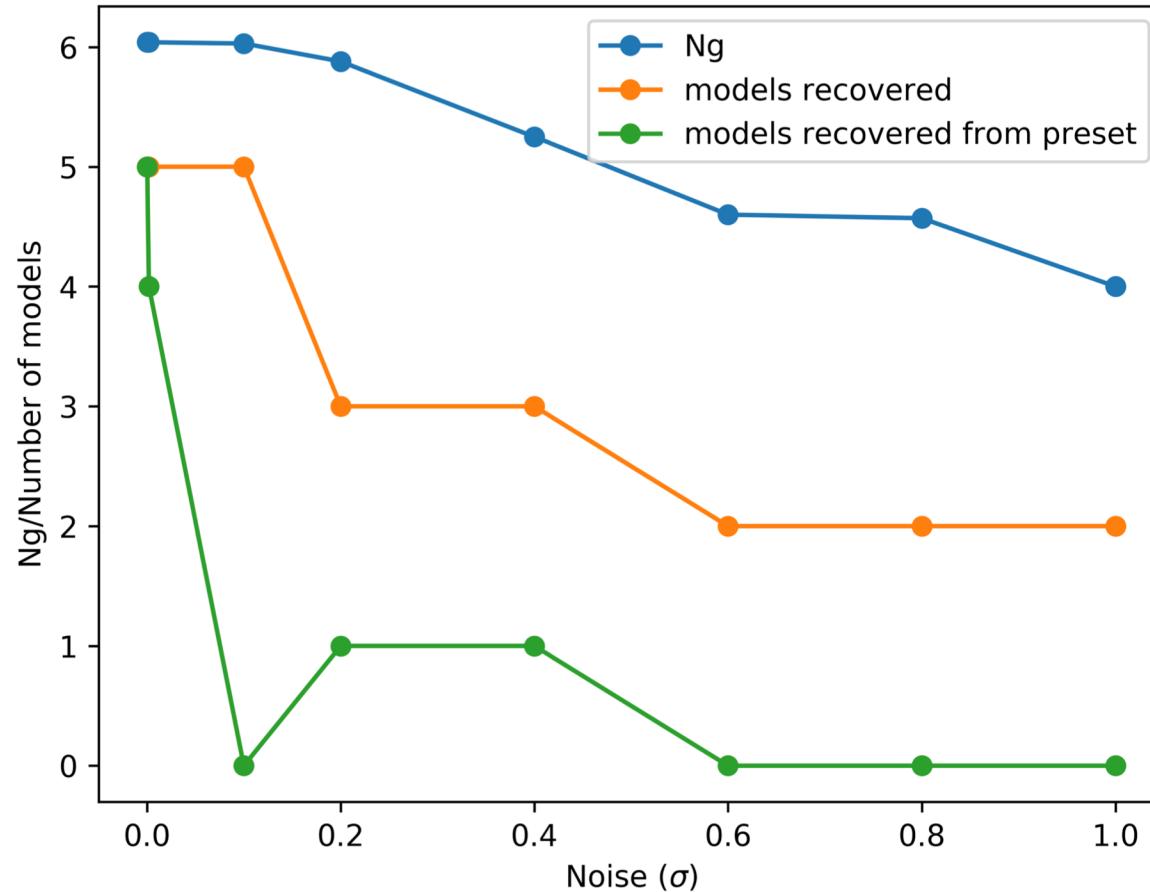


Simulated ensemble with 5 members from a library of 100 conformers

# SasView has large collection of models



# Variational Bayesian Inference on noisy data



Simulated ensemble with 5 members from a library of 100 conformers

# Information content in SAS

## From Nyqvist-Shannon sampling theorem

- The scattering curve  $I(q)$  of object with maximum diameter  $d$  is fully represented by a set of points

$$q_n = n\pi/d, \text{ where } n = (1, \dots, \infty).$$

- The number of independent data points (Shannon channels)

$$N_s = d(q_{\max} - q_{\min})/\pi$$

# Bayesian ensemble modeling: complete inference

A **model** ( $\mathbf{M}$ ) is a set a set of structures  $S$ . For a given model:

$$P(\mathbf{w}|\mathbf{D}) = \frac{P(\mathbf{D}|\mathbf{w}) * P(\mathbf{w})}{\int P(\mathbf{D}|\mathbf{w}) * P(\mathbf{w}) d\mathbf{w}}$$

The **prior** of  $\mathbf{w}$  can be selected as a flat prior or based on Boltzmann distribution

$$w_i = \frac{e^{-U(S_i)/k_B T}}{\sum_{j=1}^n e^{-U(S_j)/k_B T}}$$

The **likelihood** of  $\mathbf{w}$  can be modeled as a Gaussian distribution for each data point,  $d_i$

$$P(d_i|\mathbf{w}, \lambda) = \frac{1}{\sqrt{2\pi\varepsilon_{SAXS,i}^2}} \exp\left(-\frac{(d_i - \lambda \sum w_i S_i)^2}{2\varepsilon_{SAXS,i}^2}\right)$$

$$P(\mathbf{D}|\mathbf{w}, \lambda) = \prod_{i=1}^N P(d_i|\mathbf{w}, \lambda)$$

# Bayesian ensemble modeling

Point estimate of w's:

$$w_j^B \equiv \langle w_j \rangle_{posterior} = \int w_j P(\mathbf{w}|D) d\mathbf{w}$$

Uncertainty is qualified with a posterior divergence  $\sigma_w$  for a vector  $\mathbf{w}$ :

$$\sigma_w = \langle \Omega^2(\mathbf{w}^B, \mathbf{w}) \rangle_{posterior}$$

$$0 \leq \sigma_w \leq 1$$

Where  $\Omega^2(\mathbf{w}^B, \mathbf{w})$  is the Jensen-Shannon divergence between vector  $\mathbf{w}^B$  and  $\mathbf{w}$

$\sigma_w=0$  if there is no uncertainty in population weights

# Approximate inference with variational Bayes

$$P(\mathbf{w}|\mathbf{D}, \mathbf{S}) = \frac{P(\mathbf{D}|\mathbf{w}, \mathbf{S}) * P(\mathbf{w}|\mathbf{S})}{P(\mathbf{D}|\mathbf{S})}$$

$G(w|\alpha, S)$  is selected to make analytical derivation possible

$$G(w|\alpha, S) = \frac{\Gamma(\alpha_0)}{\sum_{i=1}^l \Gamma(i)} \prod_{i=1}^l w_i^{\alpha_i - 1}$$

This is the Dirichlet distribution where  $\sum_{i=1}^l w_i = 1$

For the prior we select an non-informative prior (Jeffrey's) Dirichlet distribution

$$P(w|S) = \frac{\Gamma(n/2)}{\sum_{i=1}^l n \Gamma(1/2)} \prod_{i=1}^l w_i^{-1/2}$$

Or informative prior based on Boltzmann distribution

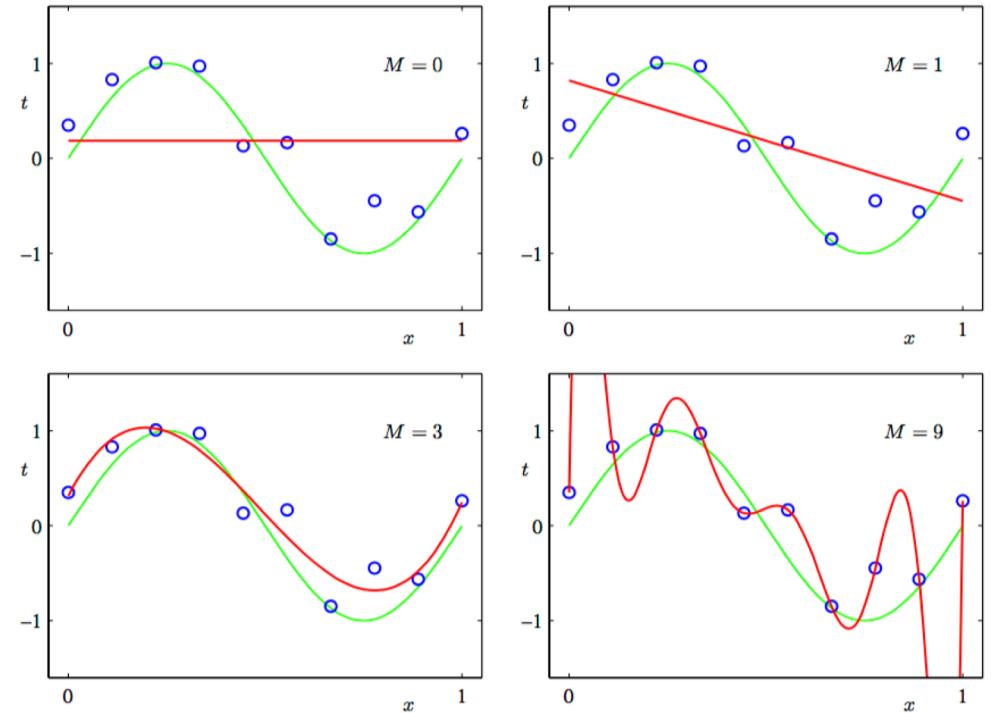
# Overfitting

More complex models will always fit the data better but risk fitting to noise

Model fitting typically involves optimization of likelihoods (Maximum likelihood)

$$\chi^2 = \frac{\sum_{i=1}^N (I(q_i)_{data} - I(q_i)_{model})^2}{\sigma_i^2}$$

Optimization based solely on  $\chi^2$  can lead to overfitting

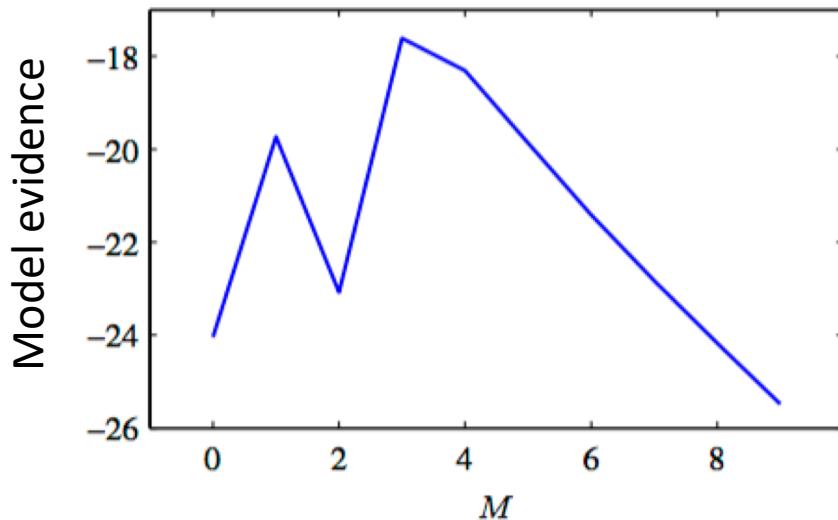


**Figure 1.4** Plots of polynomials having various orders  $M$ , shown as red curves, fitted to the data set shown in Figure 1.2.

From “Pattern recognition and machine learning” Bishop (2006)

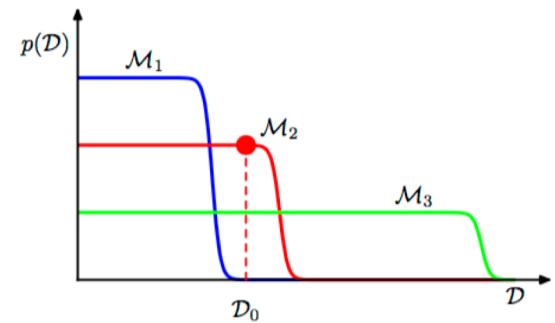
# Overfitting

Bayes to the rescue!



$$P(D|M_i) = \int P(D|\theta, M_i) * P(\theta|M_i)d\theta$$

Embodies the Occam's razor principle: simpler models are better



Figures from “Pattern recognition and machine learning” Bishop (2006)

# Calmodulin ensemble compared with PDB structures

