

# AI in SAS II

Emre H Brookes

EMBO Practical :  
Small Angle Neutron and X-ray Scattering  
from biomacromolecules in solution

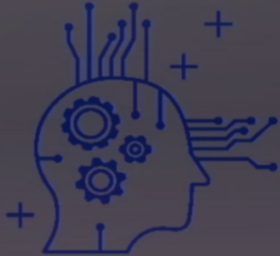
19 September 2024



# AI

## ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



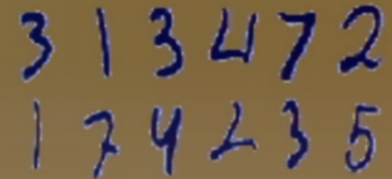
## MACHINE LEARNING

Ability to learn without explicitly being programmed

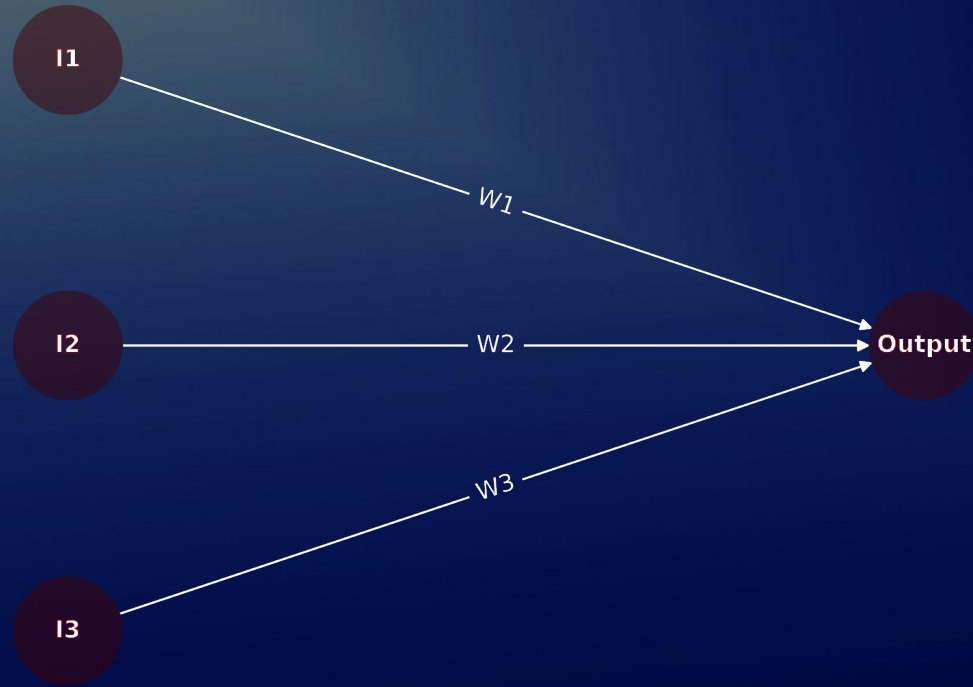


## DEEP LEARNING

Extract patterns from data using neural networks

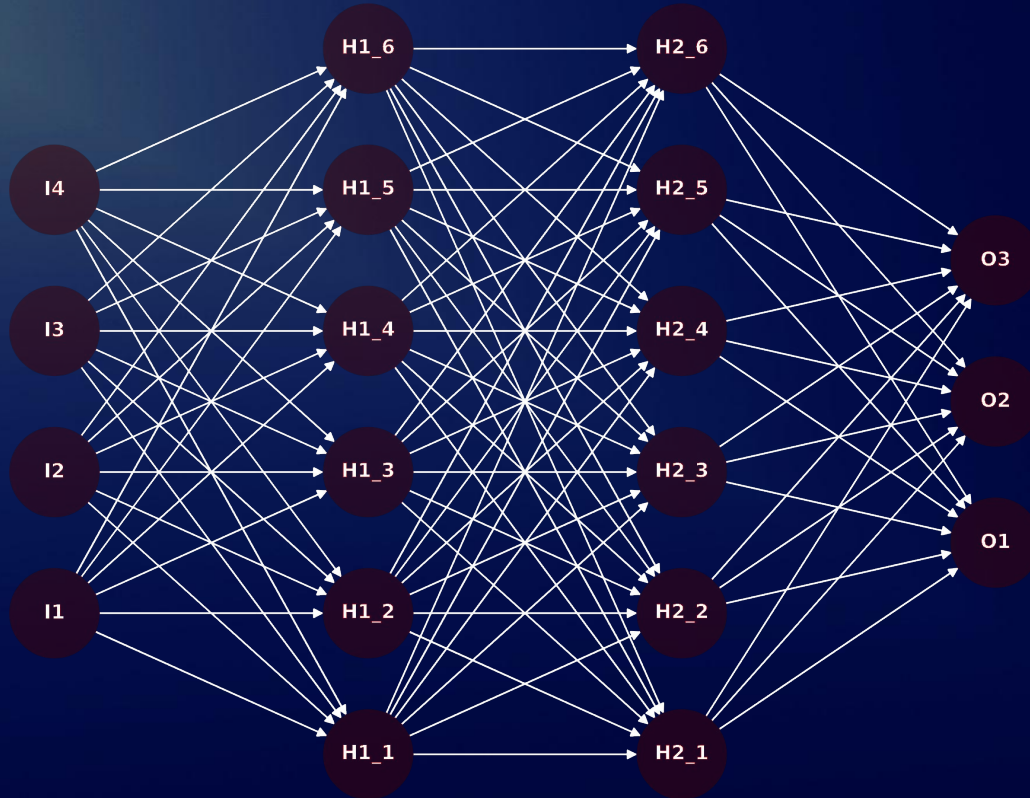


# Perceptron

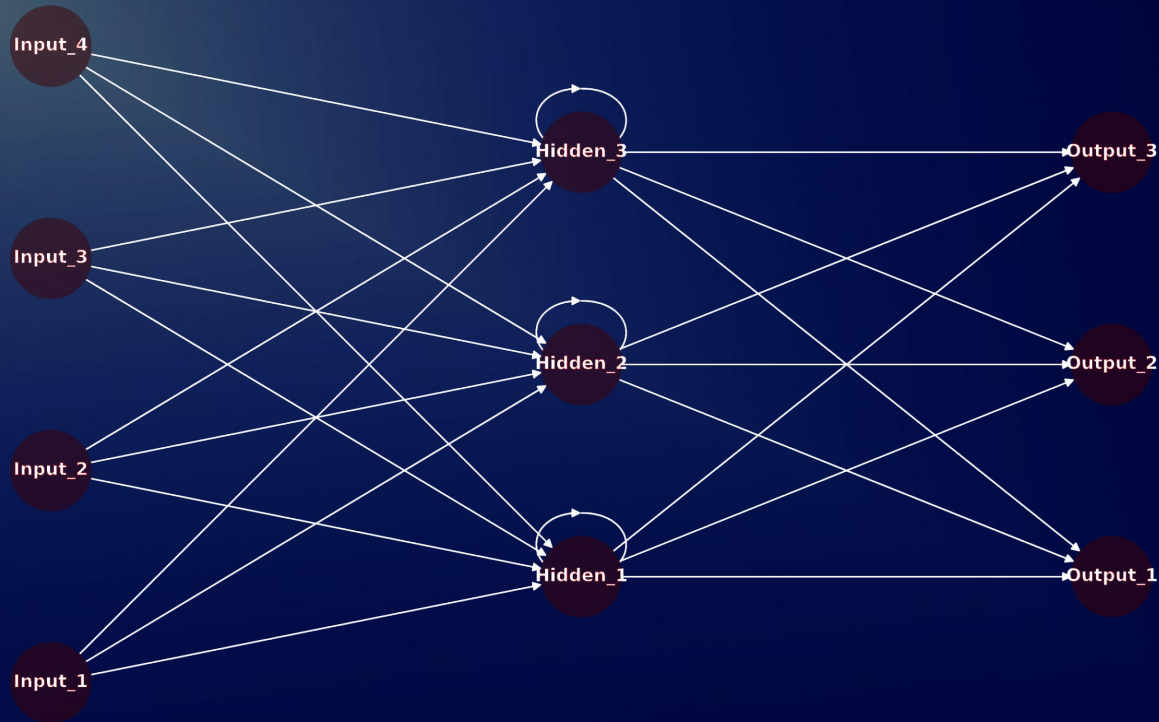


$$O = g\left(\sum_i w_i I_i\right)$$

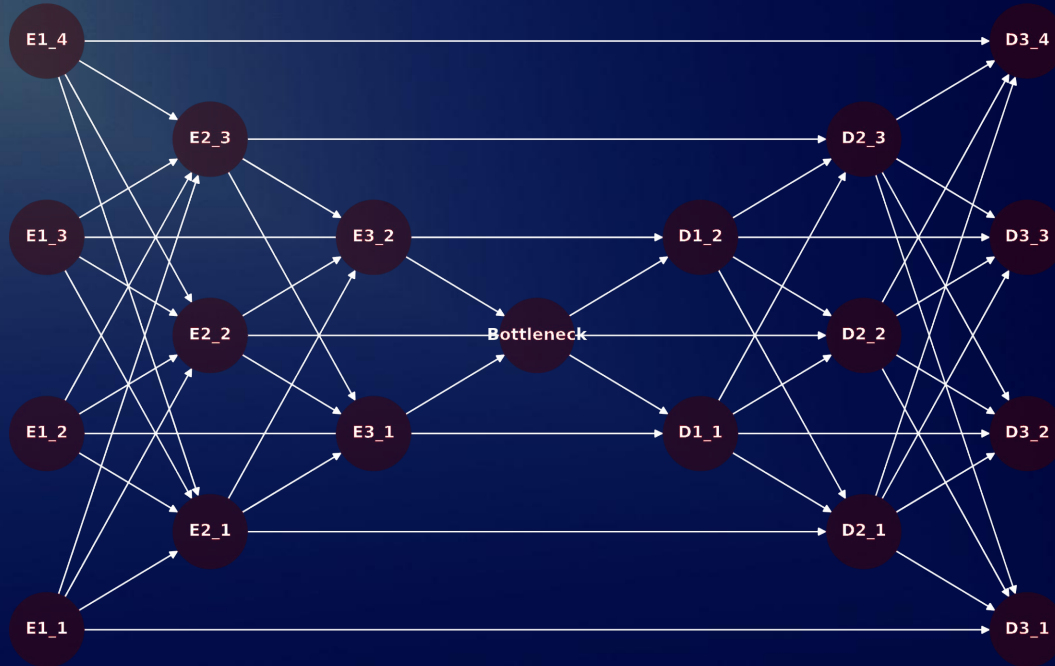
# Neural Networks - Dense 4-6-6-3



# Neural Networks - 4-3-3 Recurrent NN (RNN)

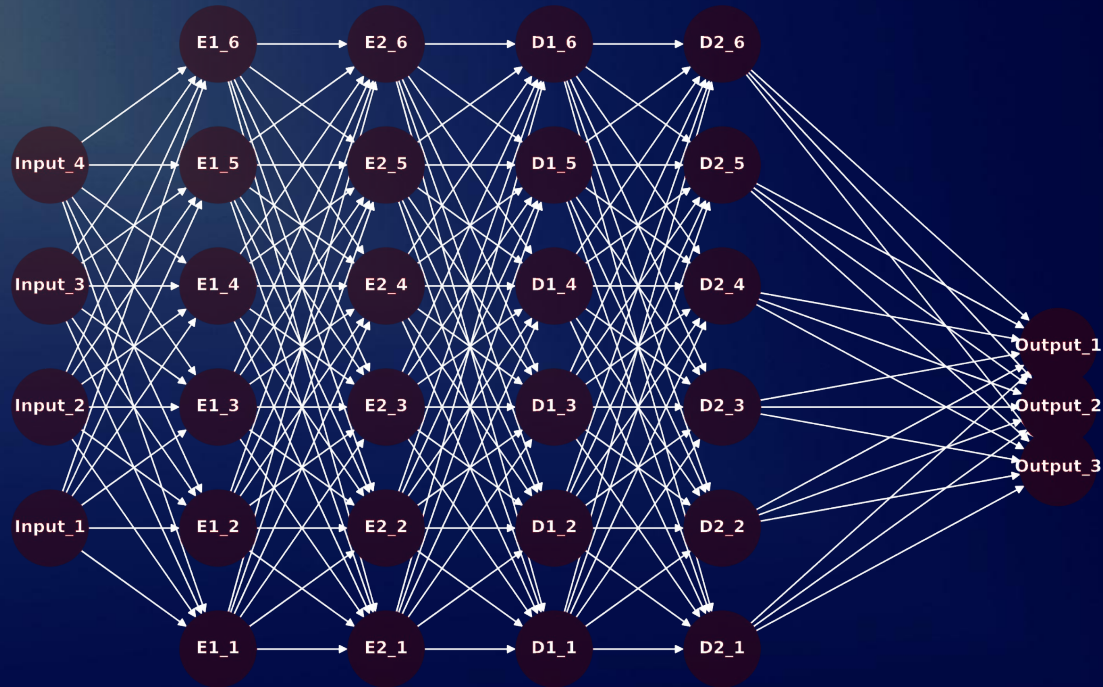


# Neural Networks - U-Net



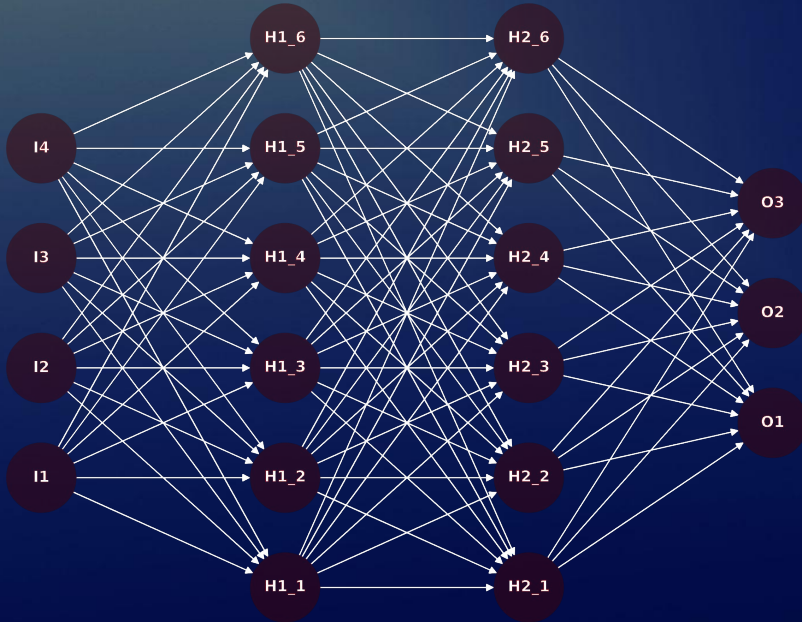
*Ronneberger et al., 2015. U-net., MICCAI 2015 proceedings, part III 18 (pp. 234-241). Springer*

# Transformer Network



*Vaswani, A., 2017. Attention is all you need. Advances in Neural Information Processing Systems.*

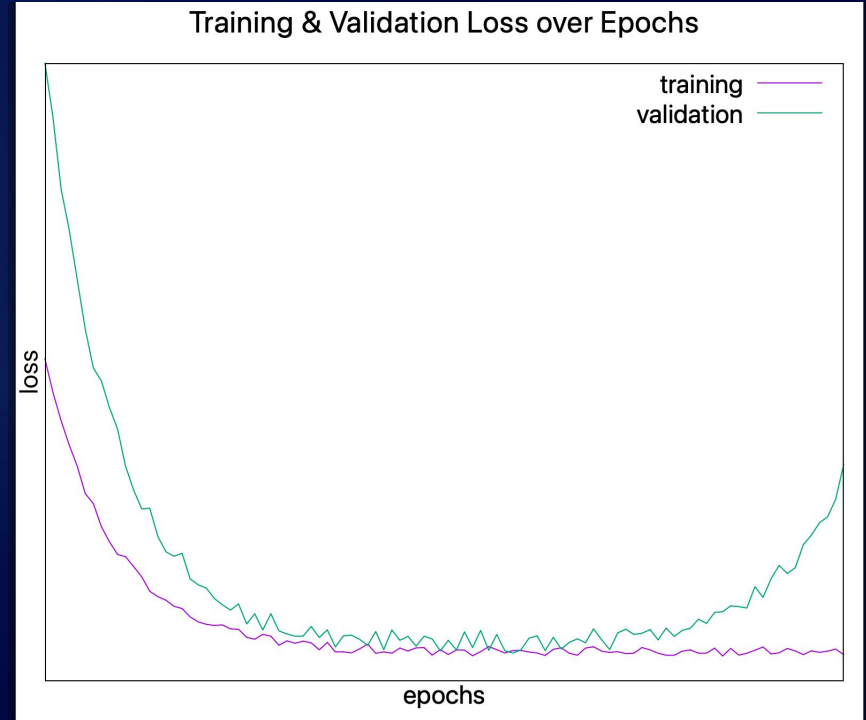
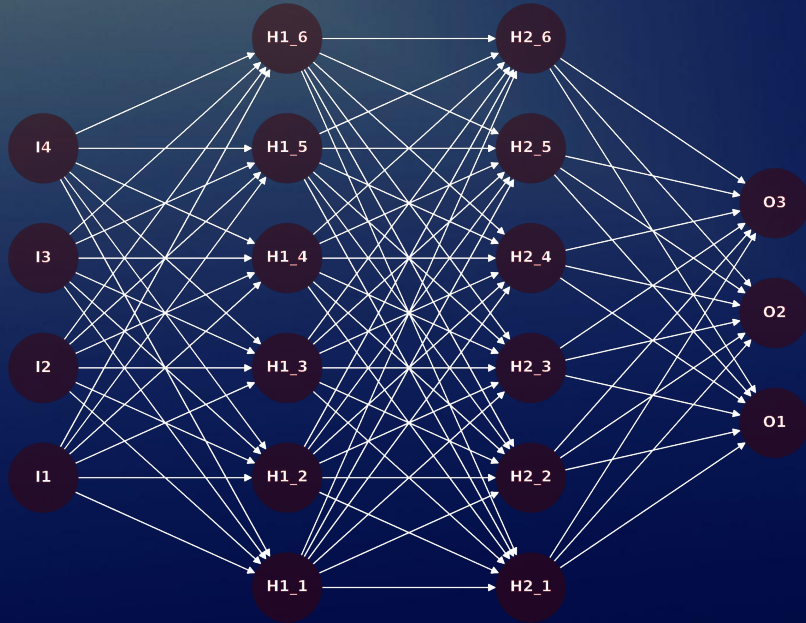
# Training



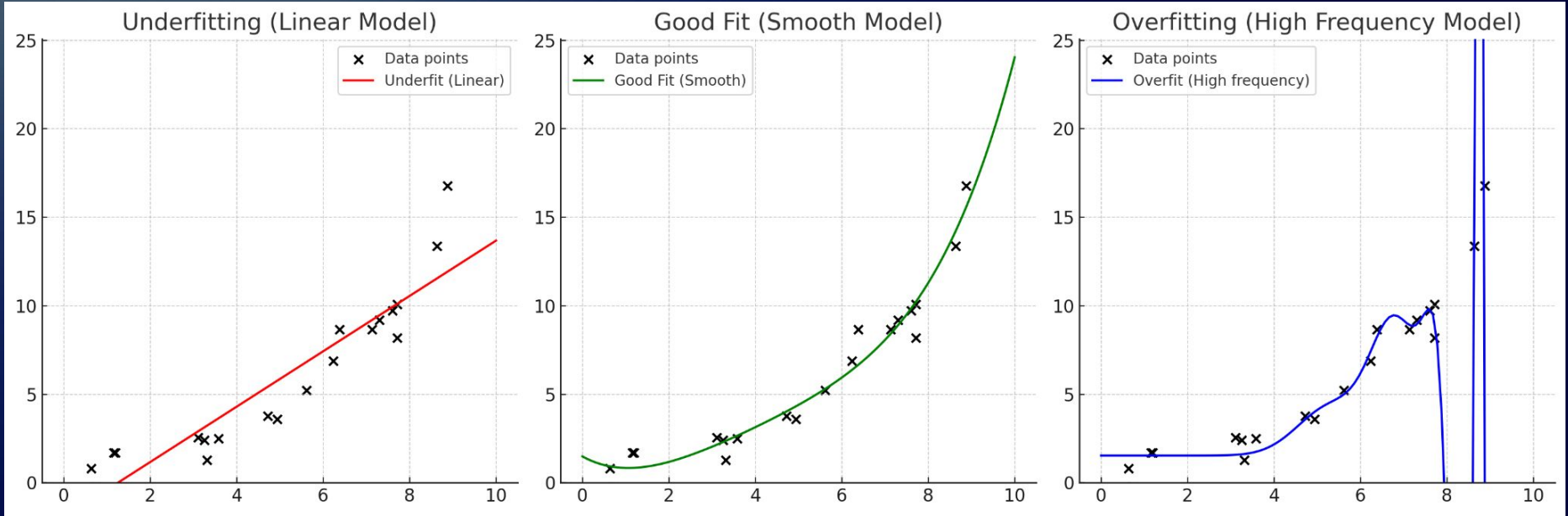
- Given:
  - training data
  - test data
- initialize random weights
- optimize the weights:
  - repeat until satisfied:
    - compute loss
    - backpropagate to update weights



# Training



# Overfitting



# ChatGPT (Generalized Pretrained Transformer)

Text



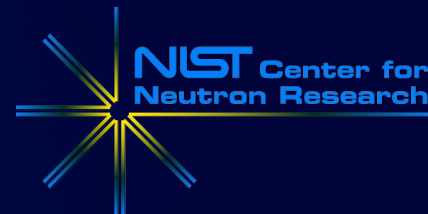
Response

- Multiple transformer blocks
- ChatGPT 3
  - *Brown, T.B., 2020. Language models are few-shot learners. arXiv preprint arXiv:2005.14165.*
  - 175 billion weights
  - 96 layers
  - The raw training set was 45TB of compressed plaintext. After filtering, it was ~570GB.
- ChatGPT 4
  - estimated 1 - 2 trillion weights

# AI/ML in Facilities

- Synchrotron optimization
  - Leemann, S.C., Liu, S., Hexemer, A., Marcus, M.A., Melton, C.N., Nishimura, H. and Sun, C., 2019. Demonstration of machine learning-based model-independent stabilization of source properties in synchrotron light sources. *Physical review letters*, 123(19), p.194801.
- Beamline optimization
  - Morris, T.W., Rakitin, M., Giles, A., Lynch, J., Walter, A.L., Nash, B., Abell, D., Moeller, P., Pogorelov, I. and Goldring, N., 2022, October. On-the-fly optimization of synchrotron beamlines using machine learning. In *Optical System Alignment, Tolerancing, and Verification XIV* (Vol. 12222, pp. 171-175). SPIE.

# NIST Center for Neutron Research



- Neutron Instrument Control Environment (NICE)
- 15 NICE instruments
- 585 total control parameters
- 818 experiment types
- Trajectory files need to be created to run the instruments
  - instrument-specific details
  - time consuming for users to master
- Given - Large amounts of documentation and example Trajectory files:

# NIST Center for Neutron Research

User text

e.g.

“Create magik trajectory angleChecks that starts with a time of 10 and loops through sampleAngle from 2 to 10 in steps of 0.25”



Properly formatted trajectory file to control the instrument

e.g.

```
{ "filePrefix": "mb111",  
  "init": [  
    ["counter.countAgainst", "'TIME'"]],  
  "loops": [{  
    "vary": [  
      ["sampleAngle", {"range": {"start": 2,  
                                "step": 0.25, "stop": 10}}],  
    ]}]
```

# NIST Center for Neutron Research

User text

e.g.

“Create magik trajectory angleChecks that starts with a time of 10 and loops through sampleAngle from 2 to 10 in steps of 0.25”



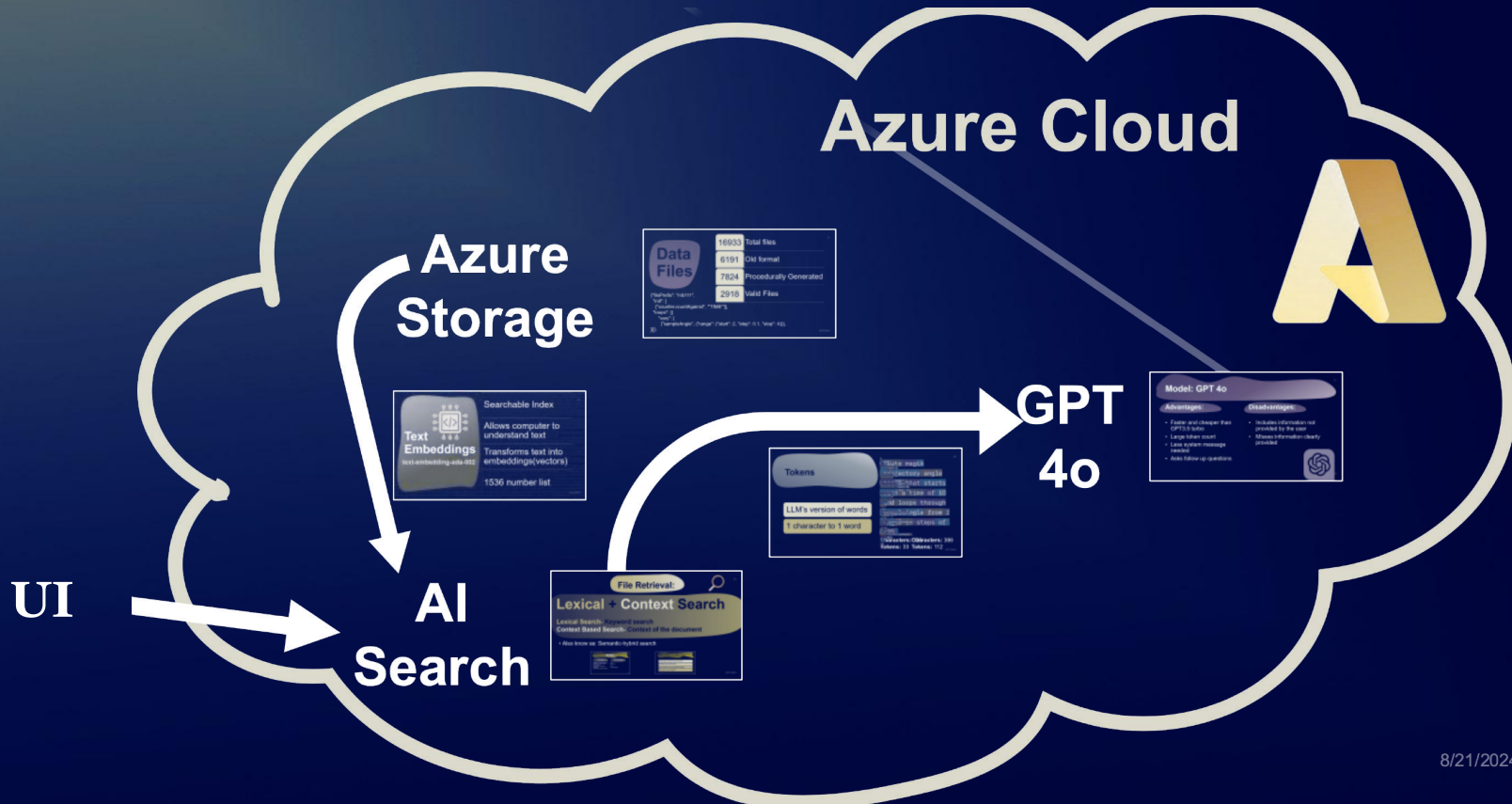
Properly formatted trajectory file to control the instrument

e.g.

```
{ "filePrefix": "mb111",  
  "init": [  
    ["counter.countAgainst", "'TIME'"]],  
  "loops": [{  
    "vary": [  
      ["sampleAngle", {"range": {"start": 2,  
        "step": 0.25, "stop": 10}}],  
    ]}]
```

*Designing and training an LLM for this would be expensive!*

# Retrieval-Augmented Generation



8/21/2024



# Retrieval-Augmented Generation

## Costs

\$6.20

To create searchable database

Per Message

\$0.02

Per Month

\$96.09

\$0.14

Per Conversation

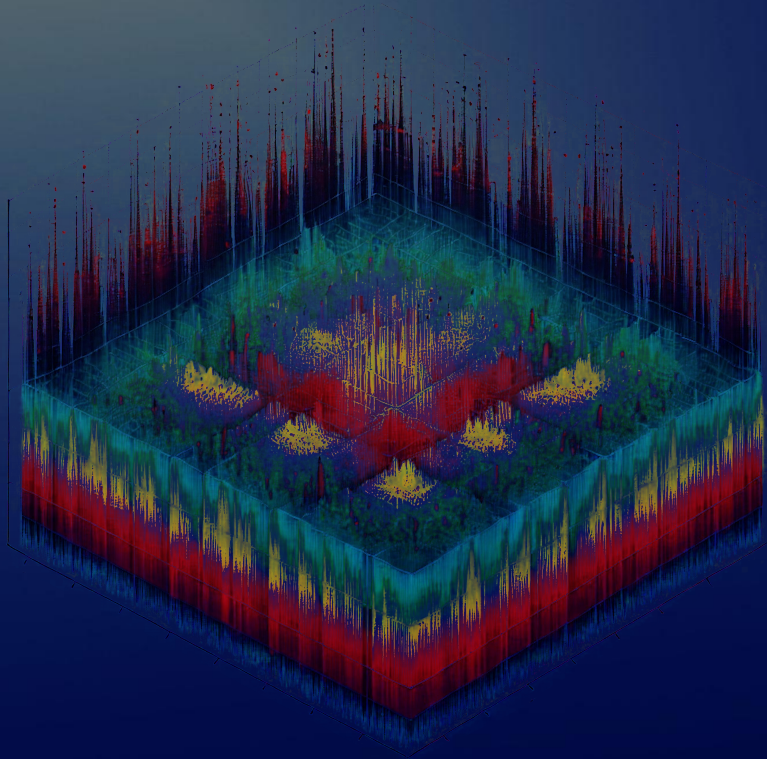
\$1,153

Per Year

- ~500 lines of code
- Data curation

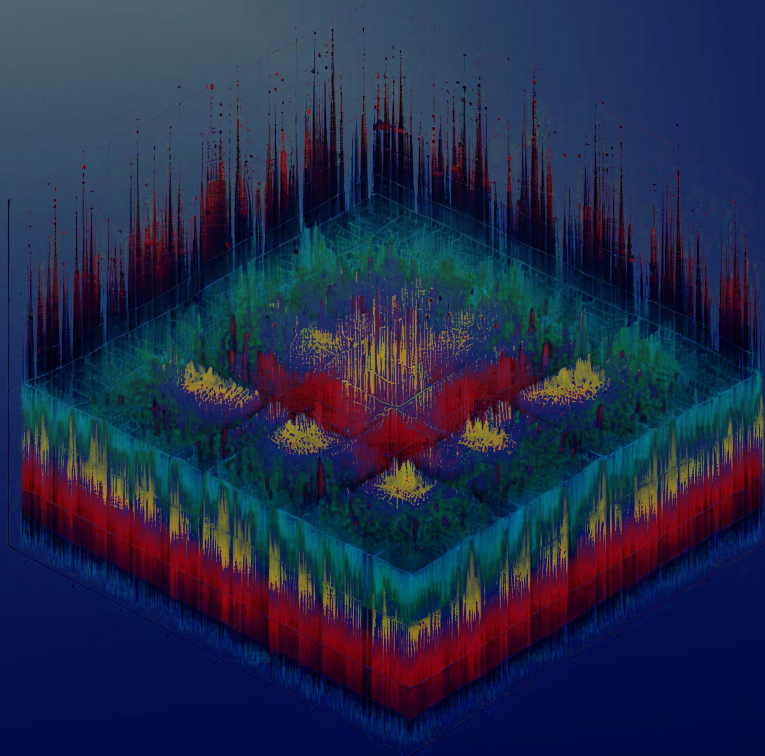


# A Common Science Problem



Relevant description of the data

# A Common Science Problem



Relevant description of the data

e.g.

- parameters
- models
- categorization

# AlphaFold / Google DeepMind

Sequence



Model

- v2 - Evoformer
- v3 - Pairformer

## Article

### Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

John Jumper<sup>1,4</sup>✉, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>✉

Nature | Vol 596 | 26 August 2021 | **583**

## Article

### Highly accurate protein structure prediction for the human proteome

<https://doi.org/10.1038/s41586-021-03828-1>

Received: 11 May 2021

Accepted: 16 July 2021

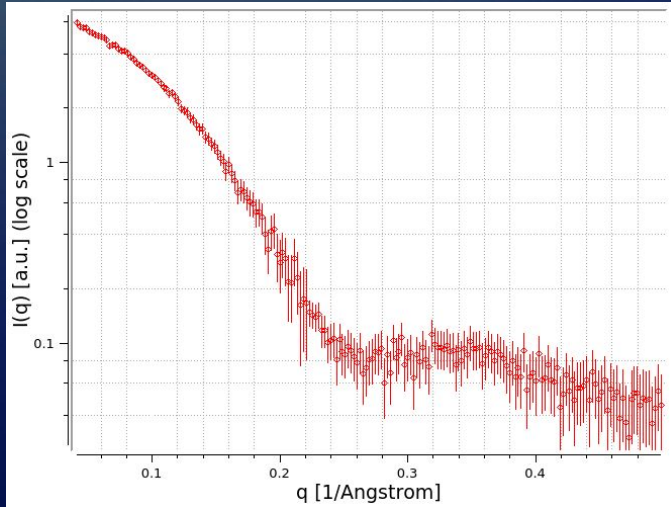
Published online: 22 July 2021

Open access

Kathryn Tunyasuvunakool<sup>1</sup>✉, Jonas Adler<sup>1</sup>, Zachary Wu<sup>1</sup>, Tim Green<sup>1</sup>, Michal Zielinski<sup>1</sup>, Augustin Židek<sup>1</sup>, Alex Bridgland<sup>1</sup>, Andrew Cowie<sup>1</sup>, Clemens Meyer<sup>1</sup>, Agata Laydon<sup>1</sup>, Sameer Velankar<sup>2</sup>, Gerard J. Kleywegt<sup>2</sup>, Alex Bateman<sup>2</sup>, Richard Evans<sup>1</sup>, Alexander Pritzel<sup>1</sup>, Michael Figurnov<sup>1</sup>, Olaf Ronneberger<sup>1</sup>, Russ Bates<sup>1</sup>, Simon A. A. Kohl<sup>1</sup>, Anna Potapenko<sup>1</sup>, Andrew J. Ballard<sup>1</sup>, Bernardino Romera-Paredes<sup>1</sup>, Stanislav Nikolov<sup>1</sup>, Rishub Jain<sup>1</sup>, Ellen Clancy<sup>1</sup>, David Reiman<sup>1</sup>, Stig Petersen<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Ewan Birney<sup>2</sup>, Pushmeet Kohli<sup>1</sup>, John Jumper<sup>1,2</sup>✉ & Demis Hassabis<sup>1,2</sup>✉

Nature | Vol 596 | 26 August 2021 | **590**

# Biological SAS $I(q)$



- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- Shape classification
- Ab initio model
- Informed model

# Biological SAS $I(q)$



- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- Shape classification
- Ab initio model
- Informed model

# Machine Learning (ML) vs “Classical”

- ML typically requires good training & test data
  - and potentially large amounts
- ML can be fast
- If the data is outside of the range of the training & test data, results will likely be wrong
- Confidence of ML results is apparently still an open question
  - Papadopoulos, G., Edwards, P.J. and Murray, A.F., 2001. Confidence estimation methods for neural networks: A practical comparison. IEEE transactions on neural networks, 12(6), pp.1278-1287.
- ML can be a “black-box”

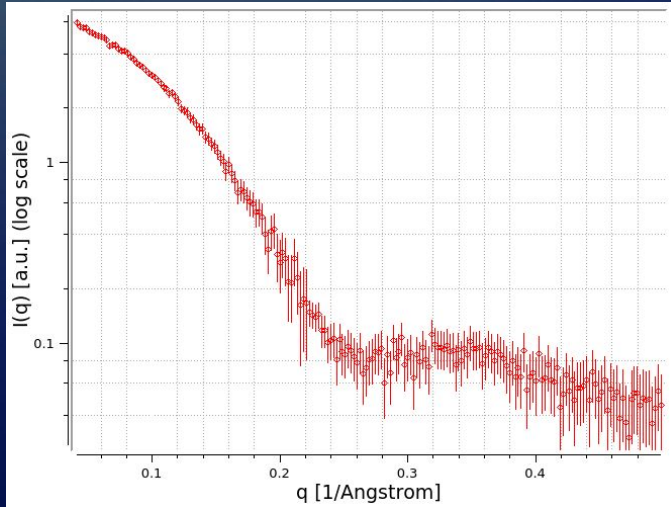
# Machine Learning vs “Classical” - ChatGPT 4o

*Tradeoffs for using machine learning vs classical approaches for modeling*

<b>Category</b>	<b>Machine Learning</b>	<b>Classical Approach</b>
<b>Data Requirements</b>	Requires large datasets	Works with limited data
<b>Interpretability</b>	Less interpretable (black-box)	Highly interpretable
<b>Generalization</b>	Good at interpolation	Good for extrapolation
<b>Computational Complexity</b>	Computationally intensive	Less computationally demanding
<b>Flexibility</b>	Highly flexible	Less flexible
<b>Model Robustness</b>	Can be less robust	More robust
<b>Development Time</b>	Longer development time	Faster development time
<b>Scalability</b>	Highly scalable	Less scalable across domains
<b>Handling Nonlinearity</b>	Great for handling nonlinearity	Limited handling of nonlinearity

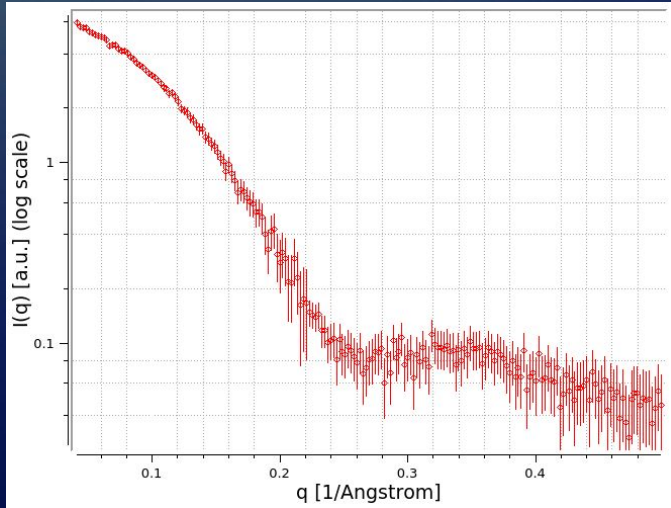


# SAS $I(q)$



- **Quality assessment**
- **$R_g$ ,  $I(0)$**
- **$D_{max}$ ,  $P(r)$**
- **Shape classification**
- **Ab initio model**
- **Informed model**

# SAS $I(q)$

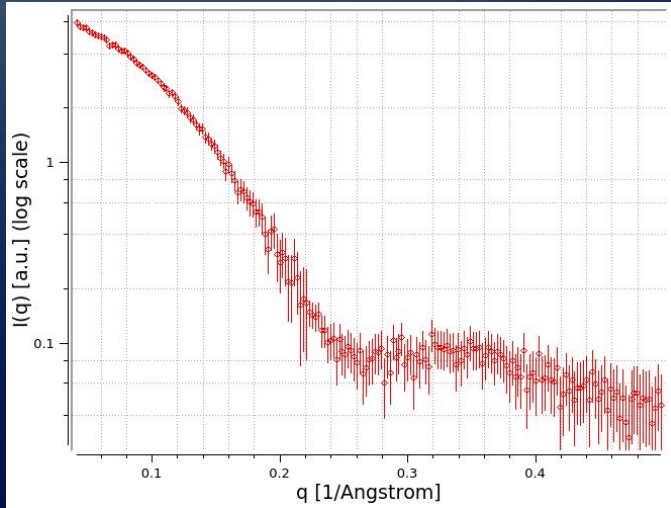


- **Quality assessment**
- **$R_g$ ,  $I(0)$**
- **$D_{max}$ ,  $P(r)$**
- **Shape classification**
- **Ab initio model**
- **Informed model**

*Prof. André Guinier*  
1911-2000  
Orsay, France

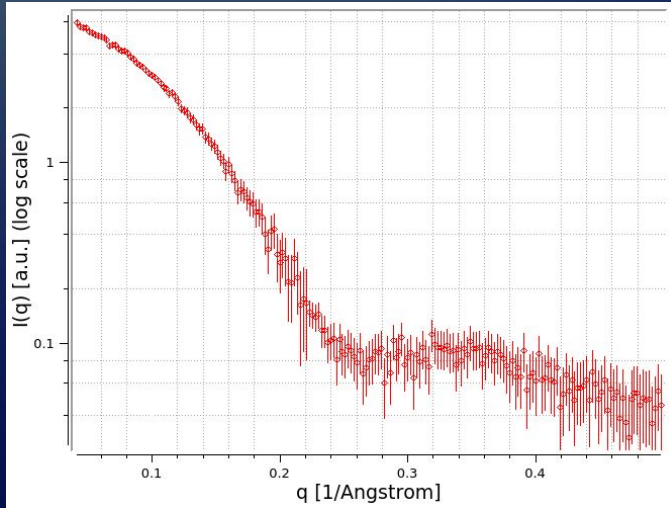


# SAS $I(q)$



- Quality assessment
- $R_g$ ,  $I(0)$
- **$D_{max}$ ,  $P(\mathbf{r})$**
- Shape classification
- Ab initio model
- Informed model

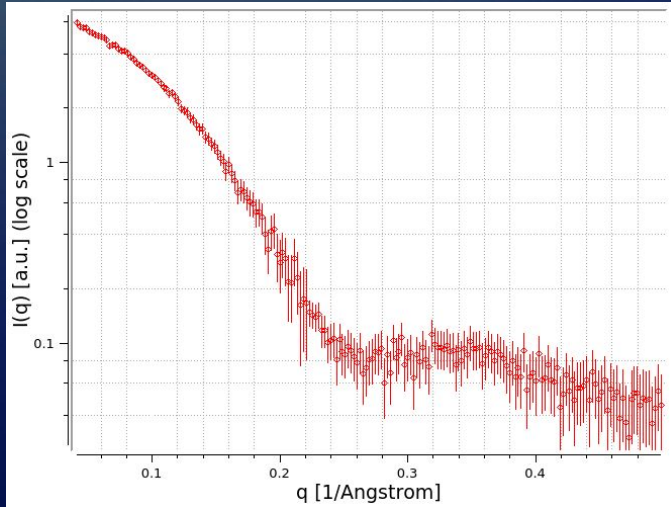
# SAS $I(q)$



- Quality assessment
- $R_g$ ,  $I(0)$
- **$D_{max}$ ,  $P(\mathbf{r})$**
- Shape classification
- Ab initio model
- Informed model

- Glatter, O. (1977) *J. Appl. Cryst.* 10, 415-421.
- GNOM - Svergun D.I. (1992) *J. Appl. Cryst.* 25, 495-503.
- Bayesian Fitting - Hansen, S. (2000) *J. Appl. Cryst.* 33, 1415-1421

# SAS $I(q)$

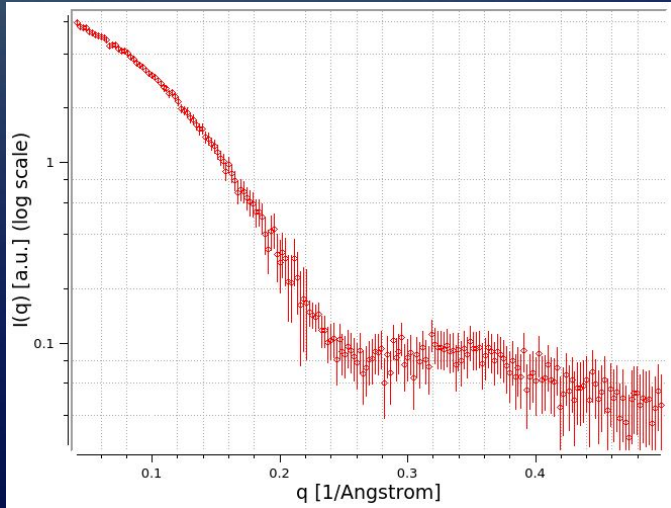


- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- **Shape classification**
- Ab initio model
- Informed model

# SAS I(q) - Shape Classification

- Classical:
  - P(r)
  - Kratky
  - Normalized Kratky
    - *Pérez, J., Vachette, P., Russo, D., Desmadril, M. and Durand, D., 2001. J. Mol. Bio., 308(4), pp.721-743.*
- ML:
  - *Franke, D., Jeffries, C.M. and Svergun, D.I., 2018. Machine learning methods for X-ray scattering data analysis from biomacromolecular solutions. Biophys. J. 114(11), pp.2485-2492.*
    - 3D feature vector space used for classification

# SAS $I(q)$

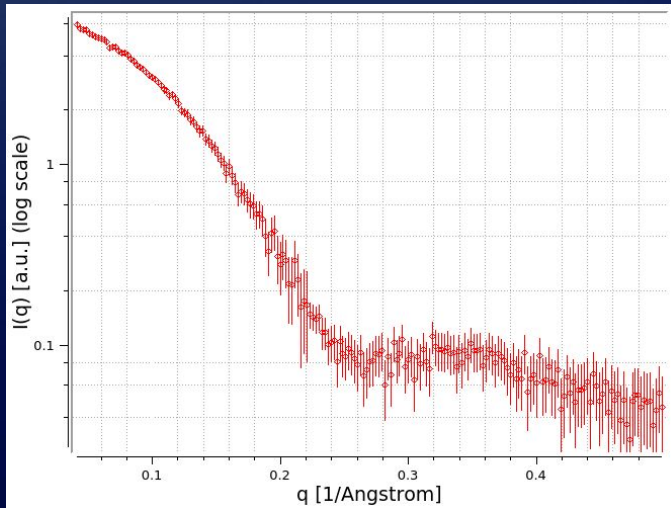


- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- Shape classification
- **Ab initio model**
- Informed model

# Information content in SAS curves

*Svergun, D.I. & Koch, M.H.J. (2003) Small-angle scattering studies of biological macromolecules in solution. Rep. Prog. Phys. 66 1735-82*

- Shannon channels =  $D_{max} \cdot q\text{-range} / \pi$
- “the number of [obtainable parameters] typically does not exceed **10–15**”

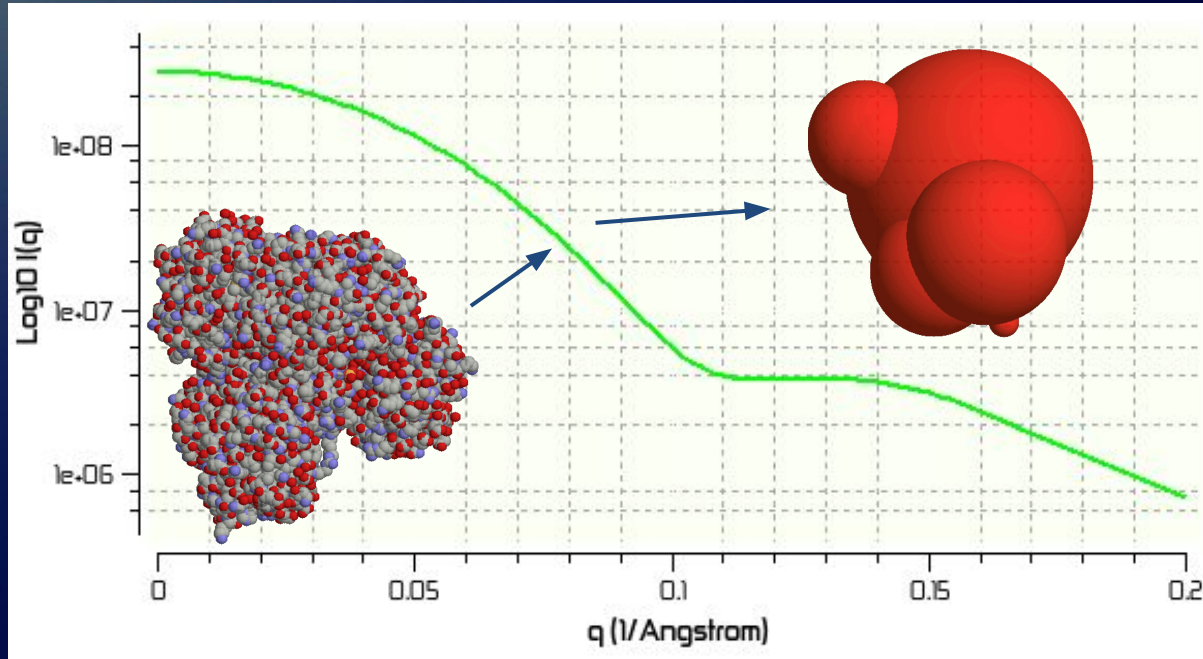


Lysozyme  $D_{max} \sim 48$  Angstroms

Shannon channels =  $48 * 0.5 / \pi \sim 8$



# Parsimonious Models



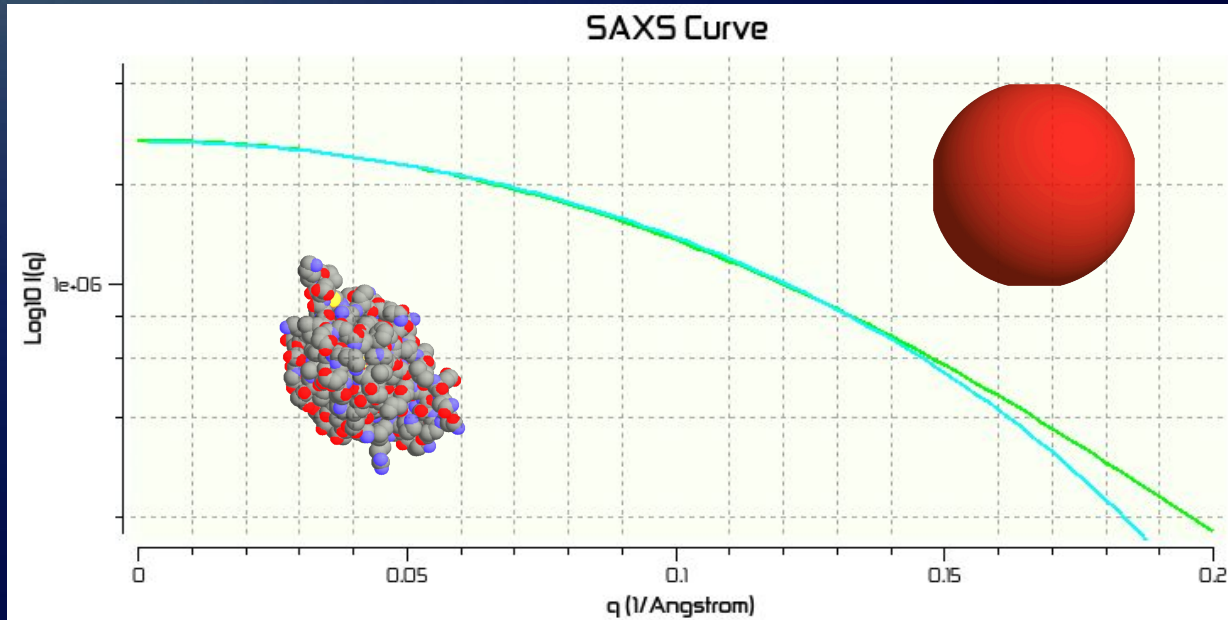
# Parsimonious Models

PDB	MW in Daltons	Description
8RAT.PDB	13,683.87	CRYSTALLOGRAPHIC STUDIES OF THE PROTEIN RIBONUCLEASE-A
1A4V.PDB	14,152.00	ALPHA-LACTALBUMIN
1DWR.PDB	17,682.20	MYOGLOBIN (HORSE HEART) WILD-TYPE COMPLEXED WITH CO
1HCO.PDB	32,279.78	HUMAN CARBONMONOXY HAEMOGLOBIN
1BEB.PDB	35,305.26	BOVINE BETA-LACTOGLOBULIN
1CTS.PDB	49,129.58	CITRATE SYNTHASE
2CGA.PDB	51,318.72	BOVINE CHYMOTRYPSINOGEN
1GZX.PDB	64,575.52	OXY T STATE HAEMOGLOBIN: OXYGEN BOUND AT ALL FOUR HAEMS
5LDH.PDB	74,917.32	ACTIVE TERNARY COMPLEX OF PIG HEART LACTATE DEHYDROGENASE WITH S-LAC-NAD
2GD1.PDB	144,427.77	OXIDOREDUCTASE(ALDEHYDE(D)-NAD(A))
1GD1.PDB	147,077.69	HOLO-GLYCERALDEHYDE-3-PHOSPHATE DEHYDROGENASE FROM BACILLUS STEAROTHERMOPHILUS
1ADO.PDB	157,287.20	FRUCTOSE 1,6-BISPHOSPHATE ALDOLASE FROM RABBIT MUSCLE
1OVA.PDB	169,965.56	UNCLEAVED OVALBUMIN

*Brookes, E., Parsimonious Spatial Models from Small Angle Scattering of Biological Macromolecules, SAS 2012, Sydney*

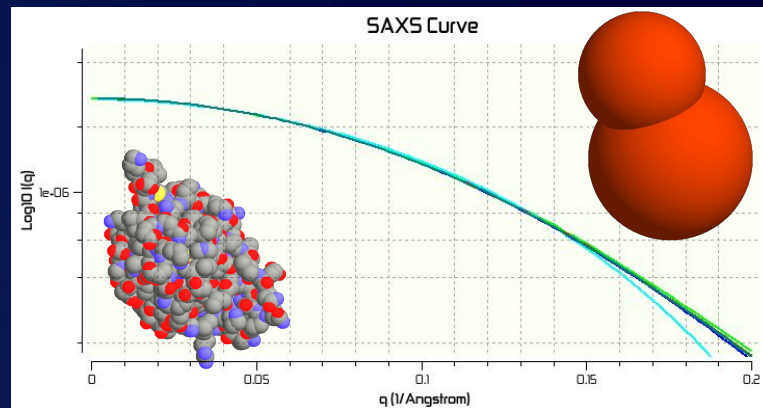
# Parsimonious Models

1A4V - 1 sphere



# Parsimonious Models

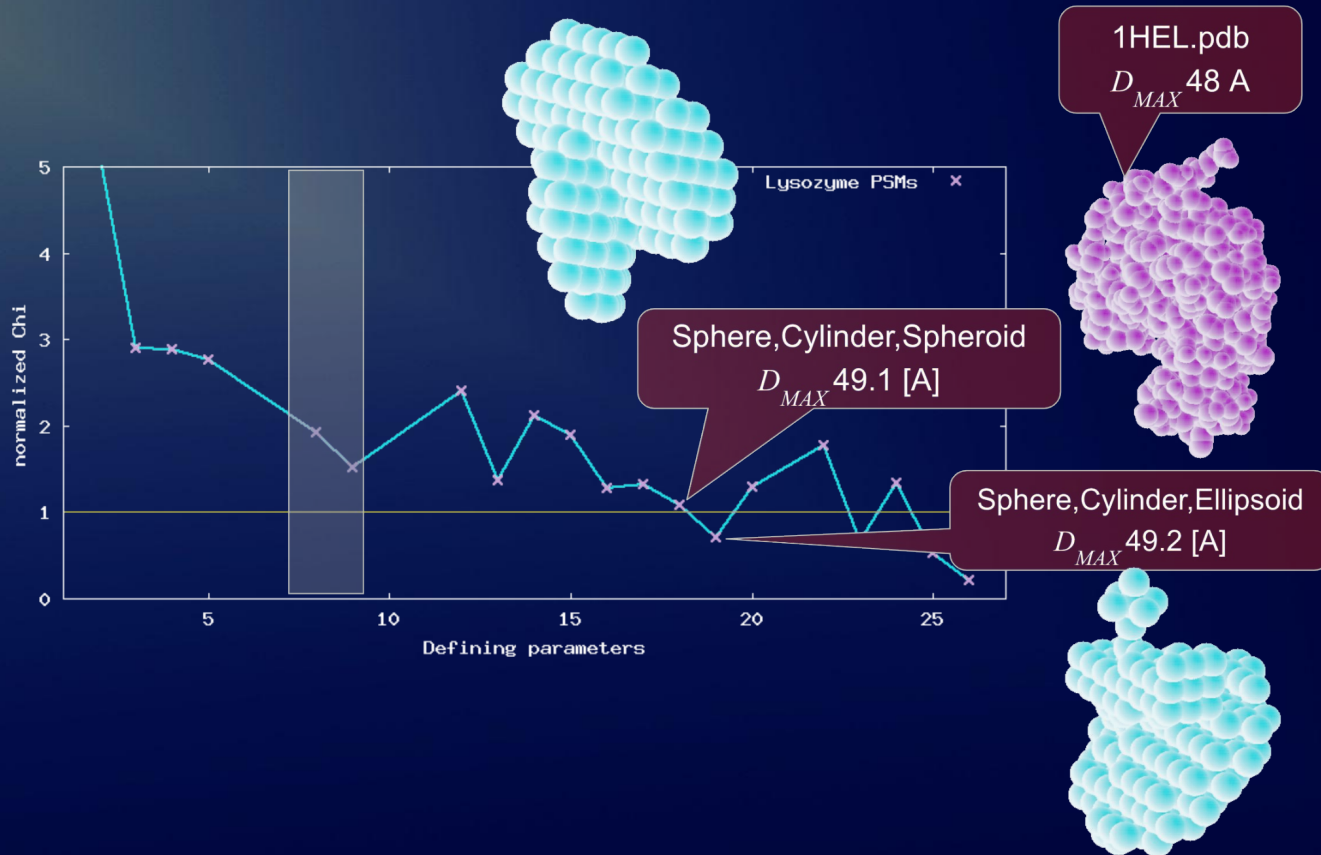
1A4V - 7 spheres



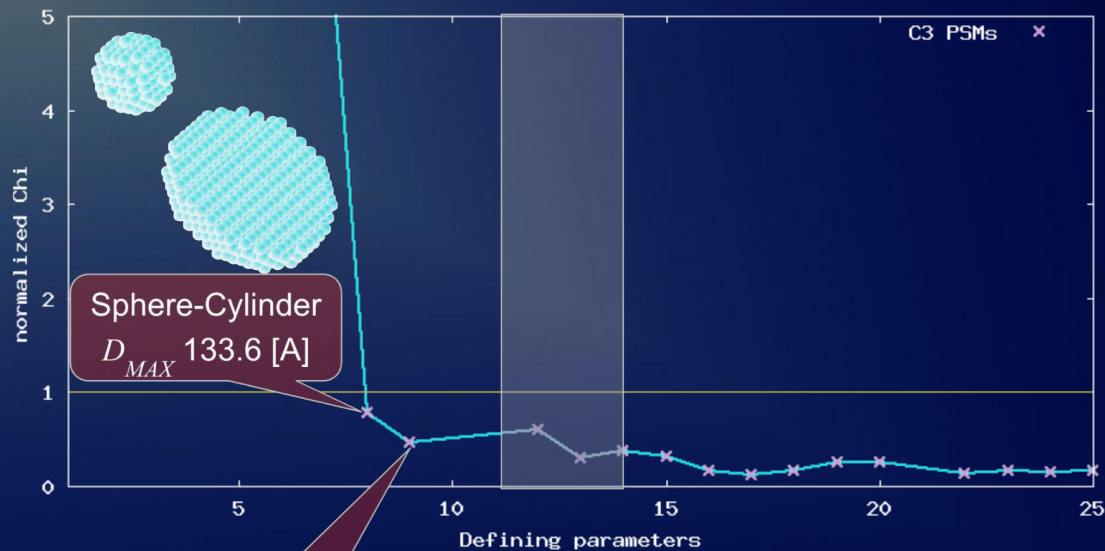
Model name	D(tr) [cm/sec <sup>2</sup> ]	Rg [nm]	Max extensions X [nm]	Y [nm]	Z [nm]	Axial ratios X:Z	X:Y	Y:Z
1A4V_1-db_1sa-10_1	1.210e-06	1.37	3.54	3.54	3.54	1.00	1.00	1.00
1A4V_1-db_2sa-10_1	1.200e-06	1.45	4.74	3.15	3.15	1.50	1.50	1.00
1A4V_1-db_3sa-10_1	1.140e-06	1.52	4.69	3.50	2.94	1.60	1.34	1.19
1A4V_1-db_4sa-10_1	1.170e-06	1.48	4.67	3.15	3.08	1.51	1.48	1.02
1A4V_1-db_5sa-10_1	1.140e-06	1.51	4.70	4.01	3.23	1.46	1.17	1.24
1A4V_1-db_6sa-10_1	1.140e-06	1.50	4.87	3.69	3.28	1.49	1.32	1.12
1A4V_1-db_7sa-10_1	1.160e-06	1.49	4.74	3.21	3.21	1.48	1.48	1.00
1A4V_1-so	1.137e-06	1.48	5.67	3.54	3.36	1.69	1.60	1.05

*Brookes, E., Parsimonious Spatial Models from Small Angle Scattering of Biological Macromolecules, SAS 2012, Sydney*

# Parsimonious Modeling



# Parsimonious Modeling

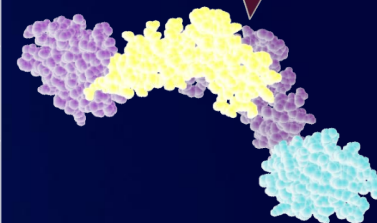


Sphere-Cylinder  
 $D_{MAX} = 133.6$  [Å]

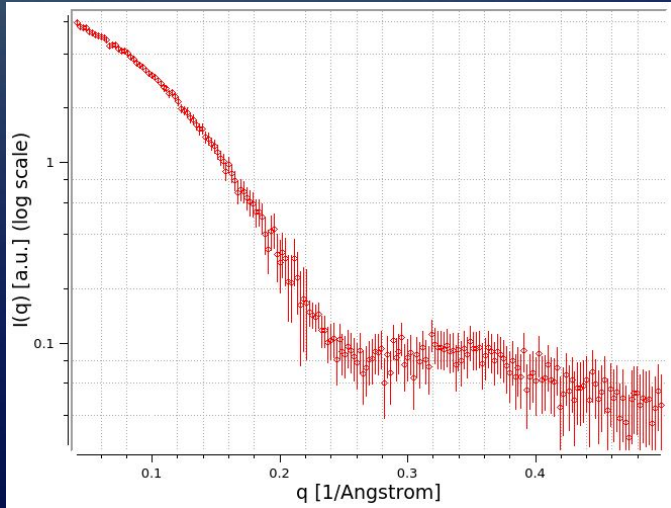
Cylinder-Cylinder  
 $D_{MAX} = 138.6$  [Å]

Homology Model

$D_{MAX} = 142$  [Å]



# SAS $I(q)$

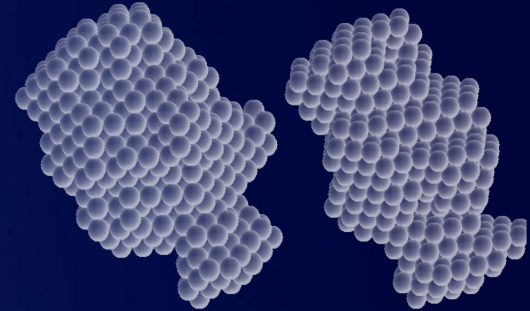
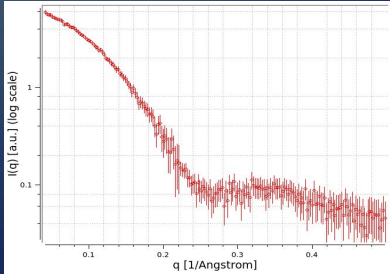


- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- Shape classification
- **Ab initio model**
- Informed model

# Ab Initio Models

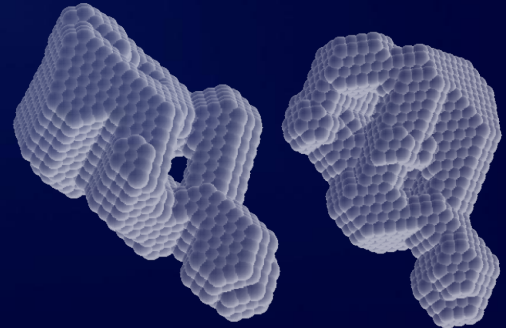
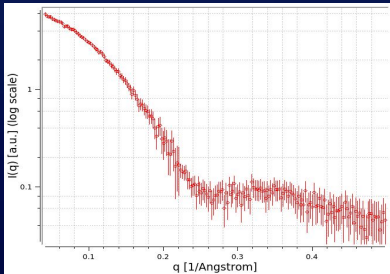
## DAMMIN

*D. I. Svergun (1999) Biophys J. 2879-2886.*



## DAMMIF

*Franke, D. and Svergun, D.I. (2009) J. Appl. Cryst., 42, 342-346.*

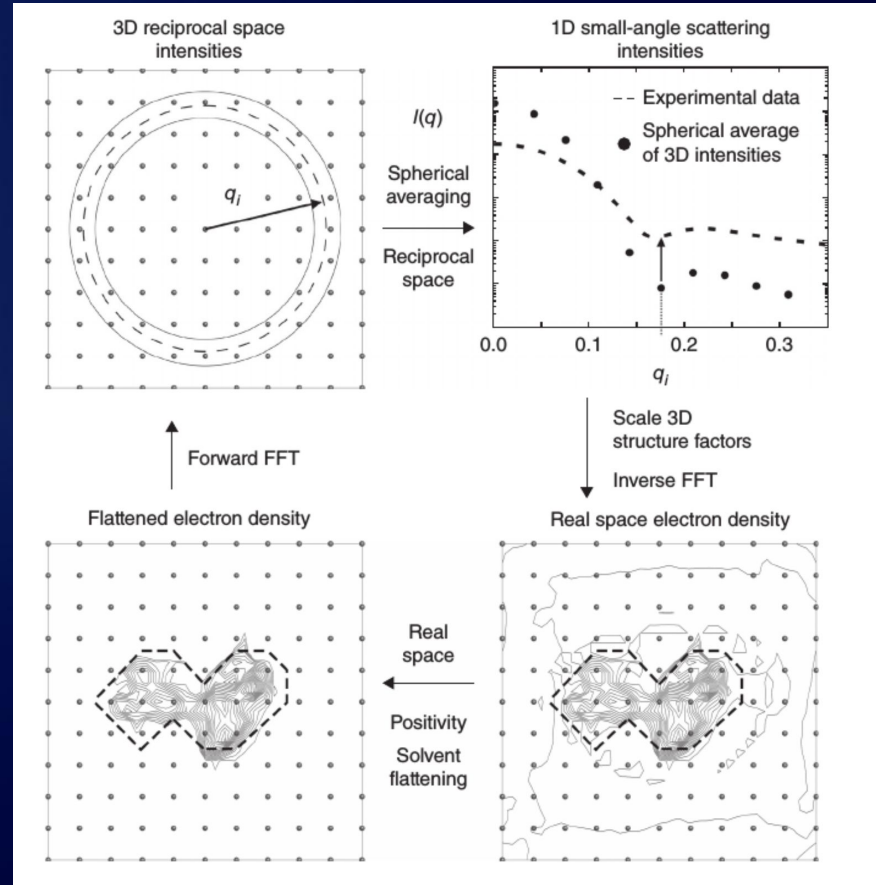




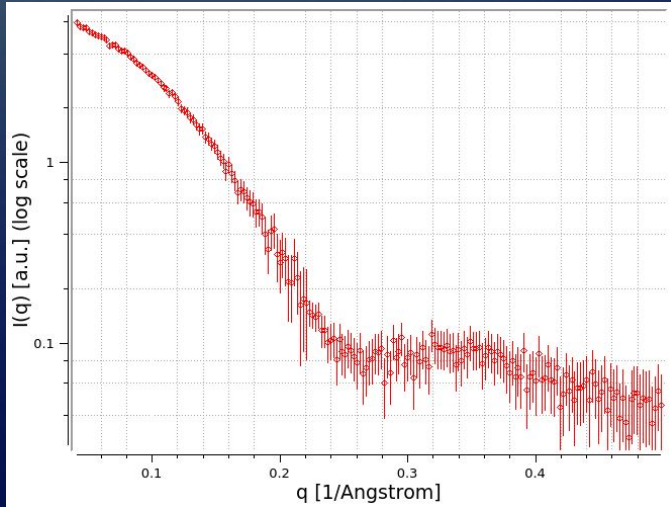
# Ab Initio Models

DENSS

*T. D. Grant (2018) Nat. Meth. 15:3 191-193*



# SAS $I(q)$

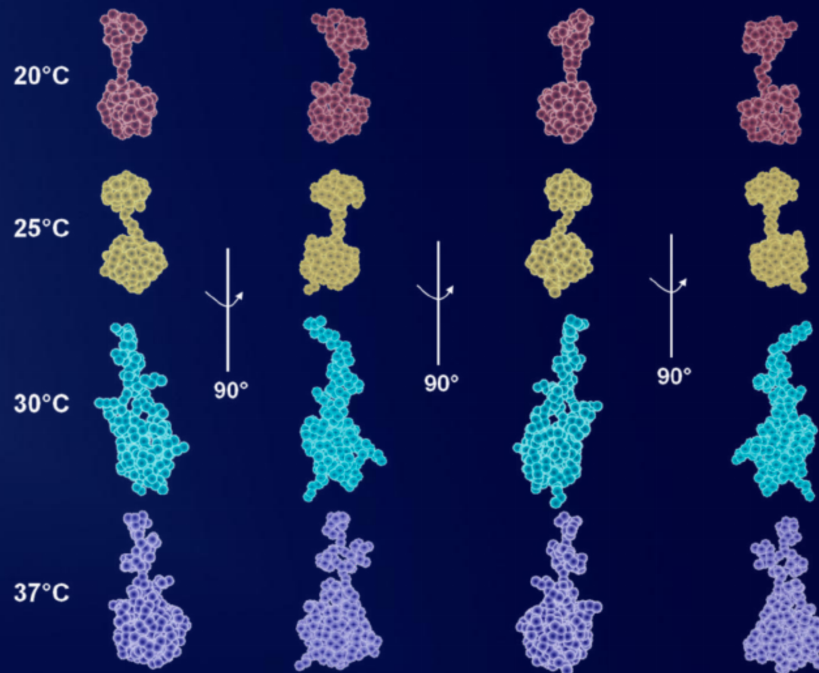


- Quality assessment
- $R_g$ ,  $I(0)$
- $D_{max}$ ,  $P(r)$
- Shape classification
- Ab initio model
- **Informed model**

# GASBOR

*Svergun, D.I. et al. (2001) Biophys. J., 80, 2946-2953.*

ab initio reconstruction of protein structure  
by a chain-like ensemble of dummy residues



Various views of the ab initio 3D models obtained using GASBOR and by averaging ten single models for each sample by using DAMAVER for rSdrFB1-4 at different temperatures.

# Information Content Revisited

*Jochen S Hub. Curr. Op. in Struct. Bio. 2018, 49:18-26*

“the interpretation of solution scattering data by computational methods is complicated by the low information content of the data, by scattering contributions from the hydration layer, and by unknown systematic errors.”

“The physical information in atomistic force fields complements the low-information SWAXS data; explicit-solvent MD may be used to predict solvent scattering, and the MD-related sampling methods may guide the structure refinement against SWAXS data.”

“Because SWAXS curves are smooth and one-dimensional (1D), they contain quite a limited amount of information. How the information is distributed over the q-range is a matter of ongoing research, but it is generally accepted that experimental SWAXS curves do not contain more than 10–30 independent data points. Hence, the number of backbone angles of biomolecules exceeds the number of independent data points of SWAXS curves by roughly two orders of magnitude. This precludes any straightforward fitting of protein structures against SWAXS data, but instead it leads to a high risk of overfitting.”

# Starting Structure + Experimental Data → Representative structure(s) that are Consistent with the Experimental Data

- Create a pool of structures from the starting structure
  - Molecular Dynamics
  - Monte Carlo
  - etc.
- Compute simulated data on each structure
- Compare the simulated data with the experimental data to choose representative structures
  - Best Fit
  - Least Squares
  - etc.

# Starting Structure + Experimental Data → Representative structure(s) that are Consistent with the Experimental Data

- Create a pool of structures from the starting structure
  - Molecular Dynamics
  - Monte Carlo
  - etc.
- Compute simulated data on each structure
- Compare the simulated data with the experimental data to choose representative structures
  - Best Fit
  - Least Squares
  - etc.

***The resulting representative structures can only be said to be consistent with the data.***

# AlphaFold / Google DeepMind

Sequence



Model

- v2 - Evoformer
- v3 - Pairformer

## Article

### Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

John Jumper<sup>1,4</sup>✉, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>✉

Nature | Vol 596 | 26 August 2021 | **583**

## Article

### Highly accurate protein structure prediction for the human proteome

<https://doi.org/10.1038/s41586-021-03828-1>

Received: 11 May 2021

Accepted: 16 July 2021

Published online: 22 July 2021

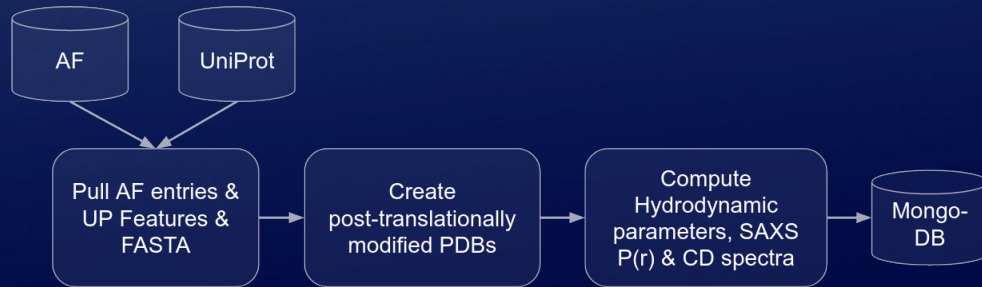
Open access

Kathryn Tunyasuvunakool<sup>1</sup>✉, Jonas Adler<sup>1</sup>, Zachary Wu<sup>1</sup>, Tim Green<sup>1</sup>, Michal Zielinski<sup>1</sup>, Augustin Židek<sup>1</sup>, Alex Bridgland<sup>1</sup>, Andrew Cowie<sup>1</sup>, Clemens Meyer<sup>1</sup>, Agata Laydon<sup>1</sup>, Sameer Velankar<sup>2</sup>, Gerard J. Kleywegt<sup>2</sup>, Alex Bateman<sup>2</sup>, Richard Evans<sup>1</sup>, Alexander Pritzel<sup>1</sup>, Michael Figurnov<sup>1</sup>, Olaf Ronneberger<sup>1</sup>, Russ Bates<sup>1</sup>, Simon A. A. Kohl<sup>1</sup>, Anna Potapenko<sup>1</sup>, Andrew J. Ballard<sup>1</sup>, Bernardino Romera-Paredes<sup>1</sup>, Stanislav Nikolov<sup>1</sup>, Rishub Jain<sup>1</sup>, Ellen Clancy<sup>1</sup>, David Reiman<sup>1</sup>, Stig Petersen<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Ewan Birney<sup>2</sup>, Pushmeet Kohli<sup>1</sup>, John Jumper<sup>1,2</sup>✉ & Demis Hassabis<sup>1,2</sup>✉

Nature | Vol 596 | 26 August 2021 | **590**

# US-SOMO AlphaFold Database

- The AlphaFold structures were predicted directly from the UniProt sequences, without any curing regarding post-translational modifications
  - Based on the UniProt annotations, removed the Initiator Methionine, Signal Peptide, and Transit Peptide(s) from the AlphaFold structures. Permuted structures with & without Propeptide(s) were also generated
- Post-translationally modified the entire AlphaFold v2 database & computed hydrodynamic, structural - incl. SAXS  $P(r)$ , & circular dichroism calculations, (~1M structures)
- <https://somo.genapp.rocks>



scientific reports View all journals Search Log in

Explore content ▼ About the journal ▼ Publish with us ▼

[nature](#) > [scientific reports](#) > [articles](#) > [article](#)

[Download PDF](#) ↓

Article | [Open access](#) | Published: 05 May 2022

## A database of calculated solution parameters for the AlphaFold predicted protein structures

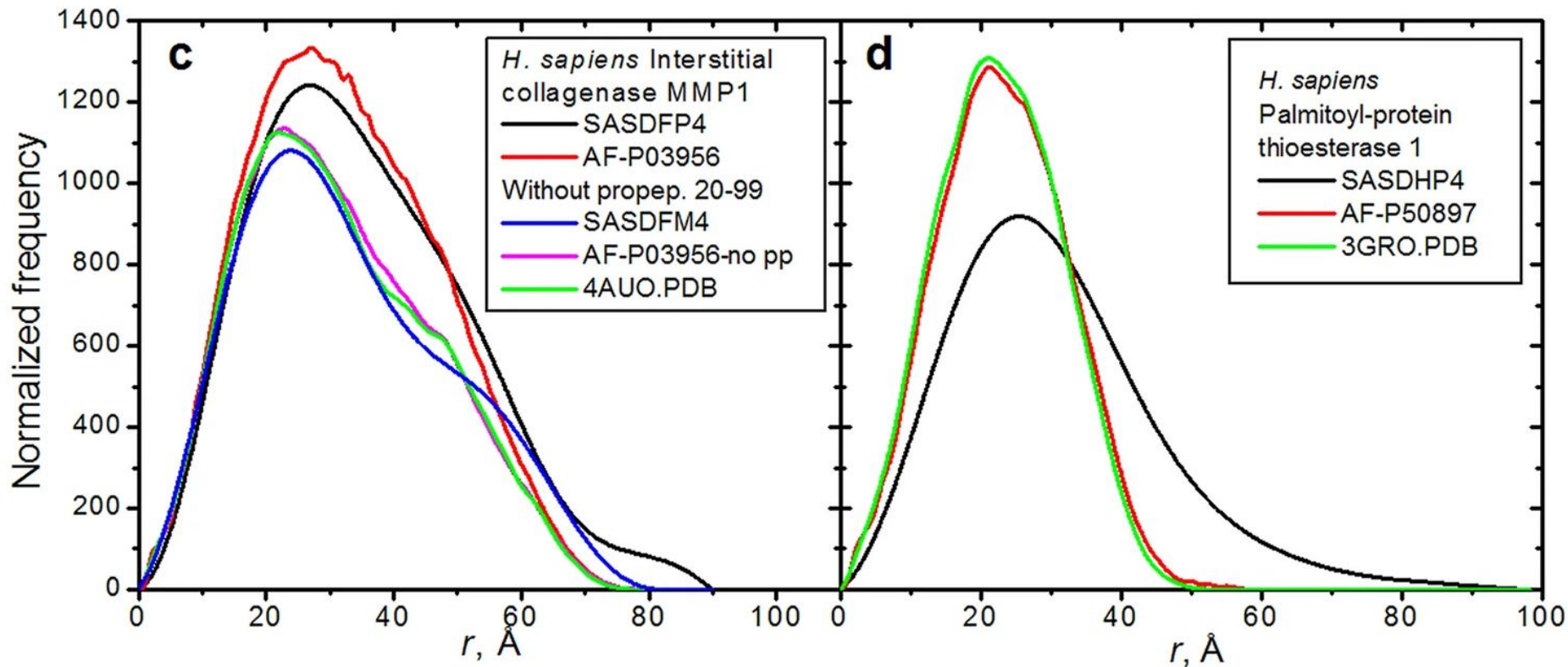
[Emre Brookes](#) ✉ & [Mattia Rocco](#)

[Scientific Reports](#) **12**, Article number: 7349 (2022) | [Cite this article](#)

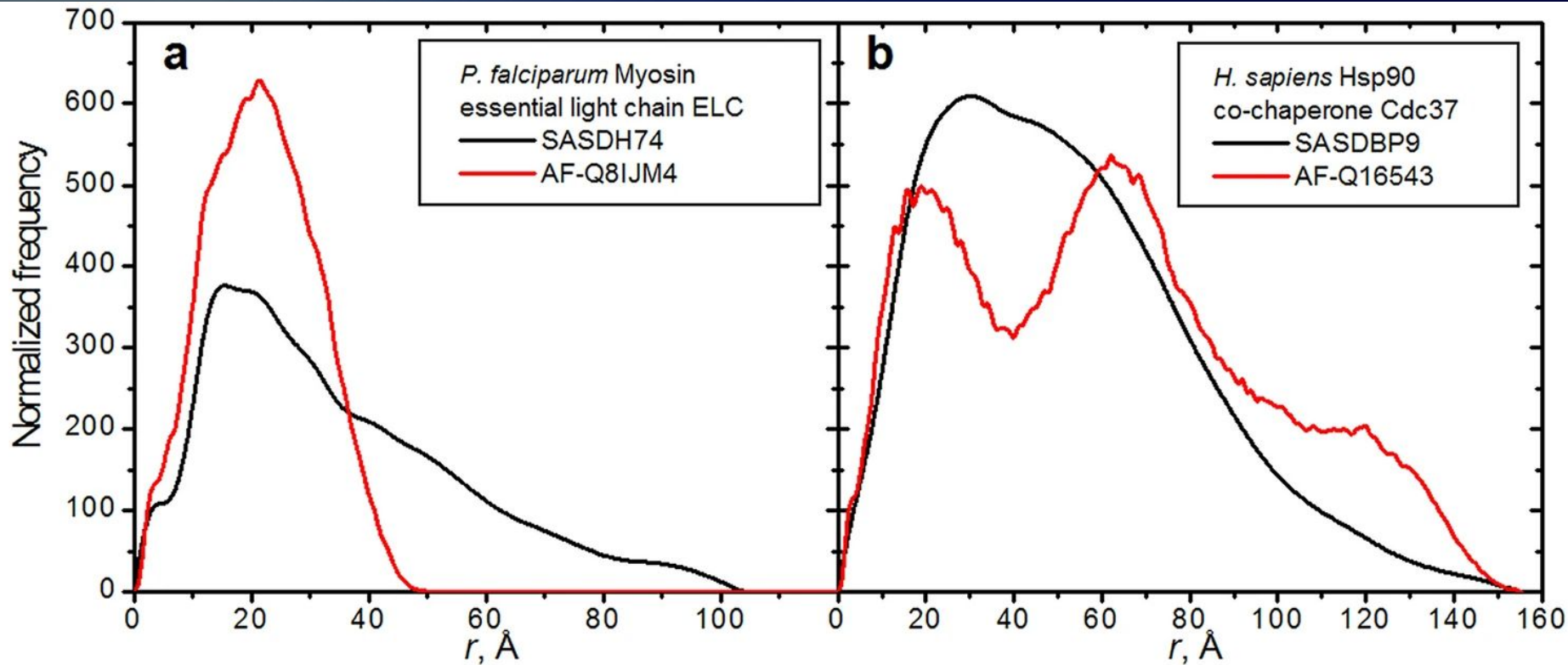
3825 Accesses | 7 Citations | 2 Altmetric | [Metrics](#)



# US-SOMO AlphaFold Database





# US-SOMO AlphaFold Database





JOURNAL OF  
APPLIED  
CRYSTALLOGRAPHY

Volume 56 | Part 4 | August 2023 | Pages 910-926  
<https://doi.org/10.1107/S1600576723005344>

OPEN  ACCESS    
Cited by  8

ISSN: 1600-5767

## 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,<sup>a\*</sup> Mattia Rocco,<sup>b</sup>  Patrice Vachette<sup>c</sup>  and Jill Trehwella<sup>d\*</sup> 

# A Fast Ensemble Modeling Method Optimizing the Fit of Protein Structures with Flexible Regions to SAXS Data

- E Brookes, M Rocco, P Vachette, J Trewella
- Inputs - Experimental data & a predicted structure, e.g. AlphaFold etc.
- Outputs - A representative ensemble
- “FAST” - NNLS fits on  $P(r)$  from MC derived pool
  - J Curtis et al. - Monomer & Complex Monte Carlo
  - Currently identifying regions of flexibility by confidence levels or user supplied
- Eventually refined in  $I(q)$  space (e.g. WAXSiS)



JOURNAL OF  
APPLIED  
CRYSTALLOGRAPHY

ISSN: 1600-5767

Volume 56 | Part 4 | August 2023 | Pages 910-926  
<https://doi.org/10.1107/S1600576723005344>

OPEN ACCESS

Cited by 8

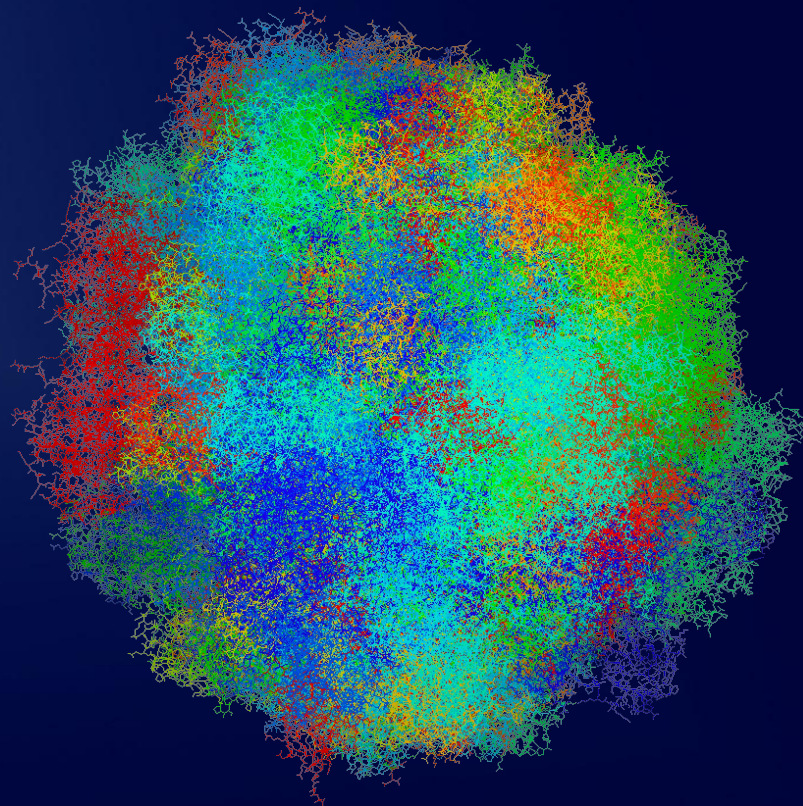
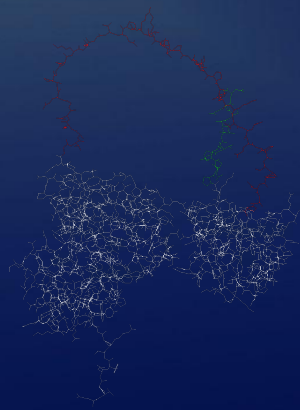
**AlphaFold-predicted protein structures and small-angle X-ray scattering: insights from an extended examination of selected data in the Small-Angle Scattering Biological Data Bank**

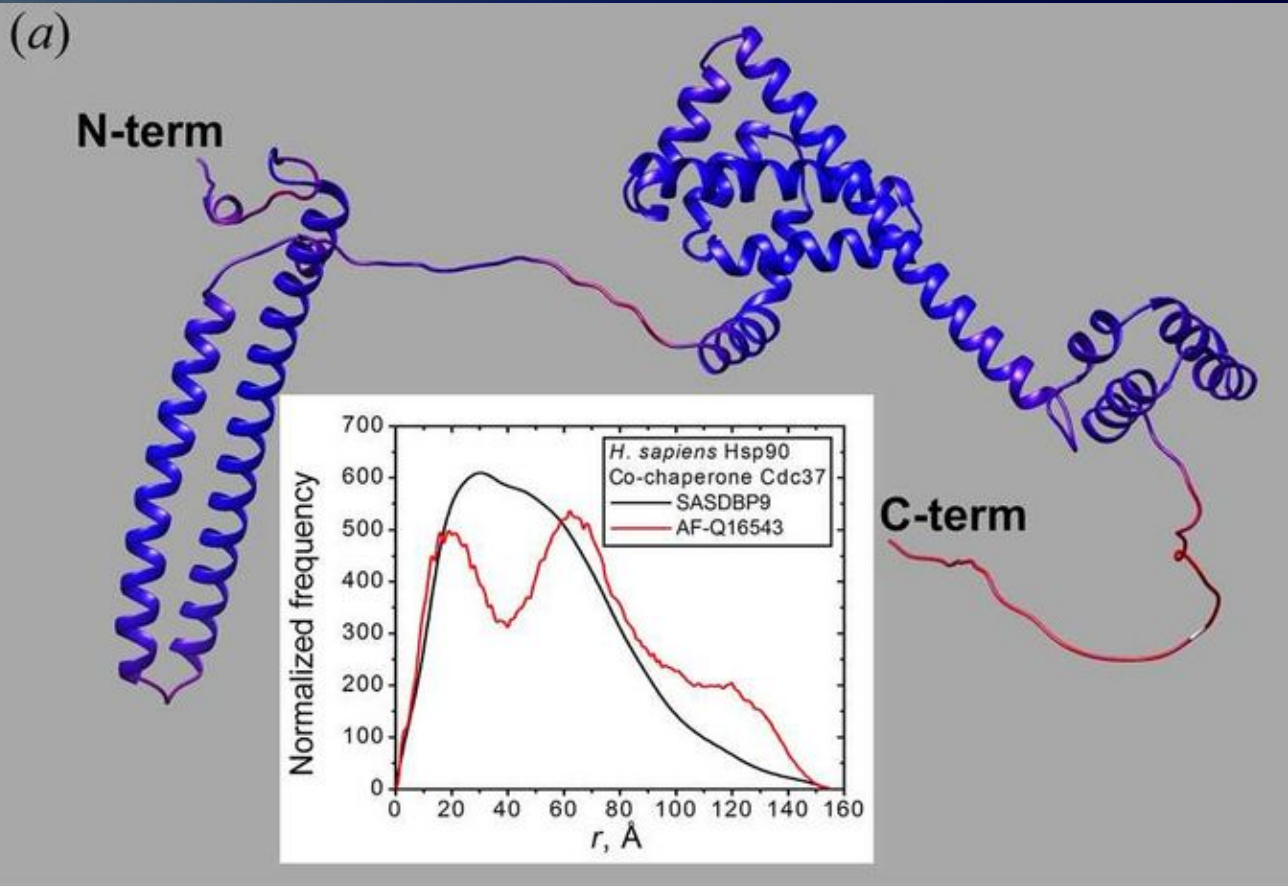
Emre Brookes,<sup>a\*</sup> Mattia Rocco,<sup>b</sup> Patrice Vachette<sup>c</sup> and Jill Trehwella<sup>d\*</sup>

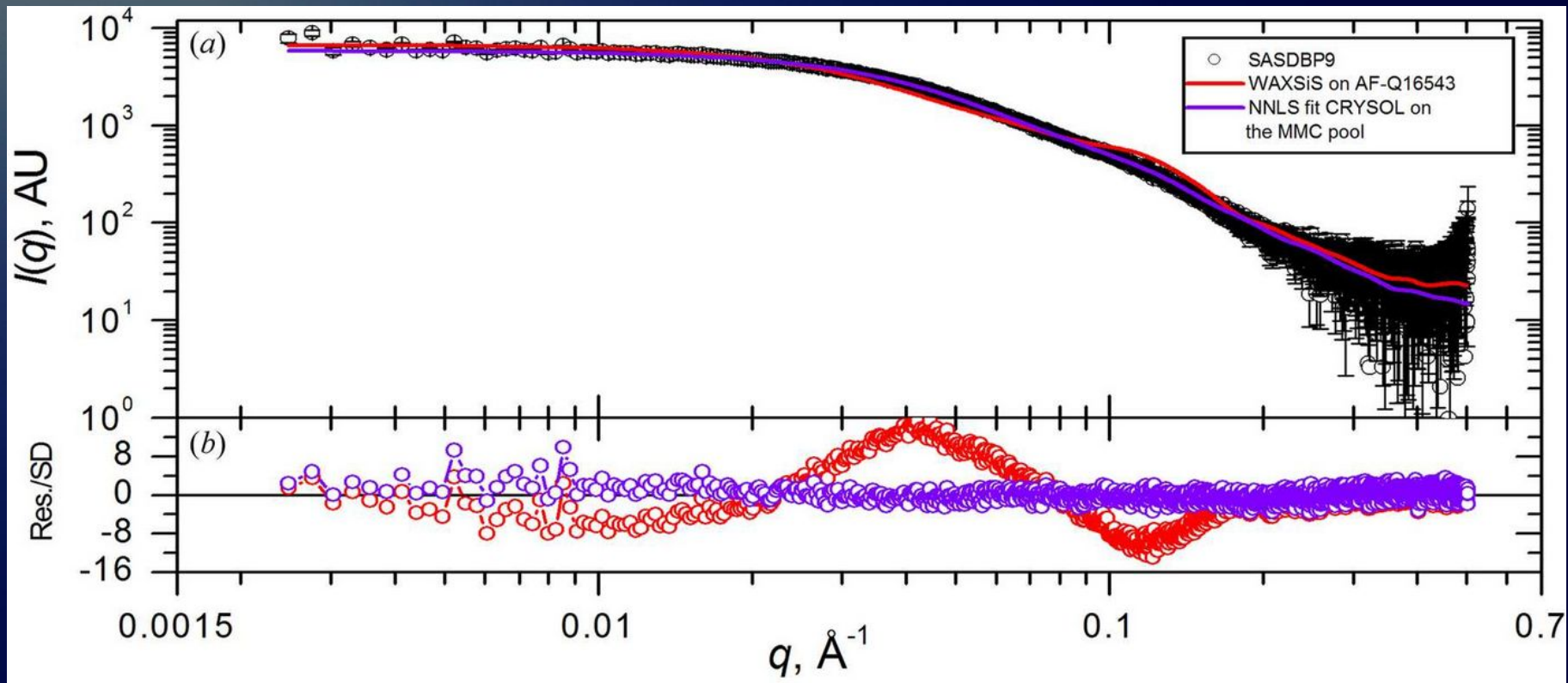
# A Fast Ensemble Modeling Method Optimizing the Fit of Protein Structures with Flexible Regions to SAXS Data

- Create a pool of structures from the starting AlphaFold (or other) structure
  - User selects flexible regions, possibly informed from AlphaFold confidence levels
  - Torsion Angle Monte Carlo
- Compute simulated data on pool of structures
  - CRY SOL
  - on selected structures - WAXSiS
    - Knight, C.J. and Hub, J.S., 2015. WAXSiS., Nucleic acids research, 43(W1), pp.W225-W230.
- Compare the simulated data with the experimental data to choose representative structures
  - Non-negatively constrained Least Squares (NNLS)

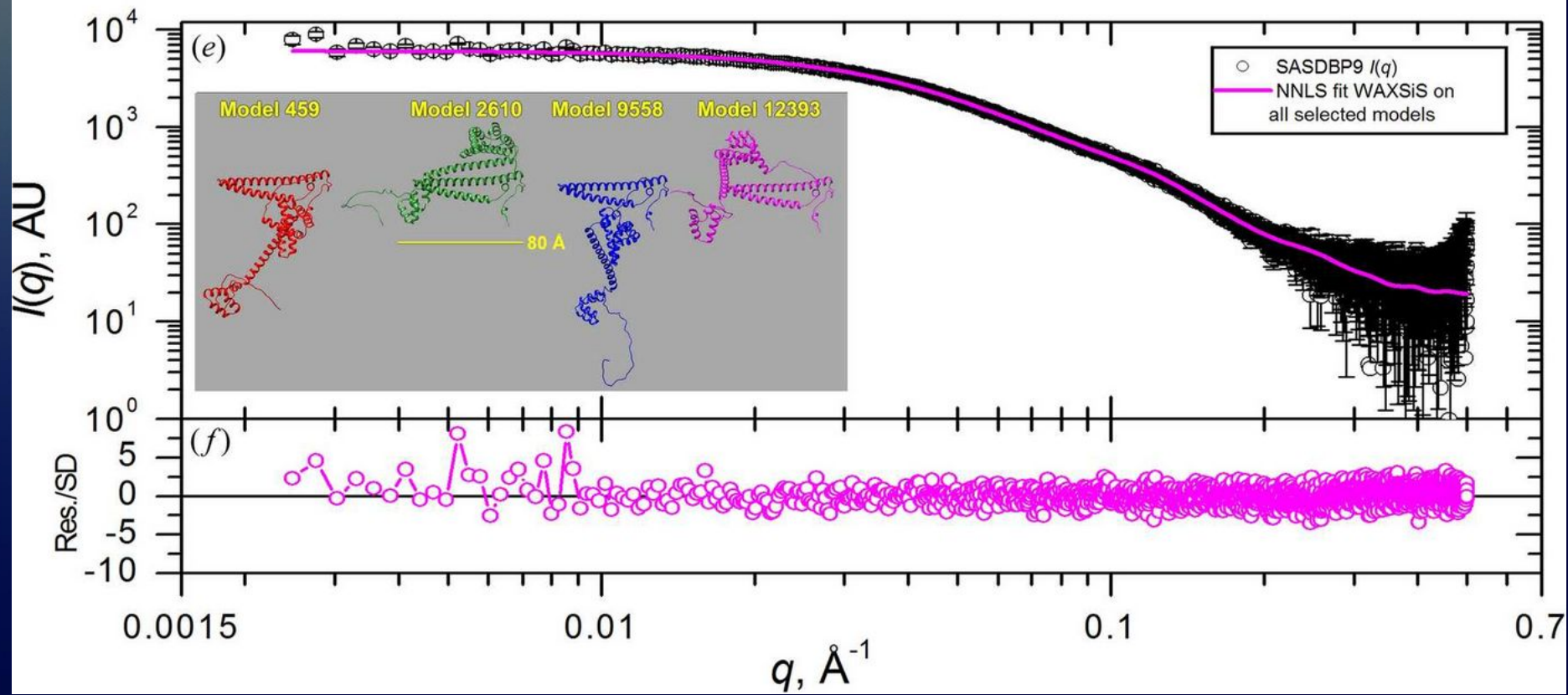
# Torsion Angle Monte Carlo on AF-Q15113











# Conclusions

- AI is a powerful tool for modeling, but is not always the best choice
  - Maslow's law of the instrument:
    - *“If the only tool you have is a hammer, you tend to see every problem as a nail.”*
- SAS I(q) data has limited information content, so, IMO, not likely a useful direct target for AI
- AI can be helpful for managing the beamline and providing user help
- AI derived structural models can provide starting structures for comparison with experimental data
- AlphaFold predicted structures do not always match solution SAS experimental data

# Acknowledgments

## Special Thanks!

- Frank Gabel

## Discussions

- Jesse Hopkins
- Robert Rambo
- Tim Snow

## Collaborators

- Joseph Curtis
- Susan Krueger
- Andreas Larsen
- Javier Pérez
- Mattia Rocco
- Jill Trehwella
- Patrice Vachette

## Infrastructure Support

- Brian Beck
- Jeremy Fischer
- Zach Graber
- Danny Havert
- Mike Lowe
- Suresh Marru
- Mark Perri
- George Turner

## Funding

- NIH K25GM090154, GM120600
- NSF CHE-1265817, OAC-1740097, OAC-1912444
- NSF XSEDE TG-MCB140255, TG-MCB170057
- NSF ACCESS MCB170057



SASBDB

Small Angle Scattering Biological Data Bank

XSEDE

Extreme Science and Engineering  
Discovery Environment

Jetstream2

EPSRC

Engineering and Physical Sciences  
Research Council



*Funded by a grant from the National  
Institute of General Medical Sciences of the  
National Institutes of Health*

ACCESS SGX3 CCP-SAS

Thank you for listening!

*Questions?*

# Publications

- *Wright DW, Nan R, Hui G, Curtis JE, Brookes EH, Perkins SJ. CCP-SAS - Novel Approaches for the Atomistic Modeling of Small Angle Scattering Data in Biology. Biophysical journal. 2015 January 27; 108(2):191a.*
- *Brookes E, Rocco M. A database of calculated solution parameters for the AlphaFold predicted protein structures. Scientific reports. 2022 May 5;12(1):7349.*
- *Brookes EH, Rocco M. Beyond the US-SOMO-AF database: a new website for hydrodynamic, structural, and circular dichroism calculations on user-supplied structures. European Biophysics Journal. 2023 Jul;52(4):225-32.*
- *Brookes, E., Rocco, M., Vachette, P. and Trewhella, J., 2023. 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. Journal of Applied Crystallography, 56(4), pp.910-926.*