Fitting challenges

David Veesler

University of Washington Department of Biochemistry <u>dveesler@uw.edu</u> <u>https://twitter.com/veeslerlab</u>

Outline

[Methods for backbone tracing into moderate resolution structures. What are the problems, how are they approached, what are the solutions? How do we validate the methods and the results?]

Single particle cryoEM at nearatomic resolution



TRPV1 3.4 Å Liao et al, Nature 2013

T20S

2.8 Å



Computational Method Single-particle cryo-EM X-ray crystallography Å 3Å 4 Å 5 Å

Main goal of structure determination



MYPRQTEINSEQVENCE



Goals of structure determination

- Determining the position of all the atoms in a structure
- Model refinement aims to:
 - Identify and correct errors
 - Improve fit to data
 - Improve model geometry
- Each atom is described by four parameters:
 - cartesian coordinates (x,y,z) & temperature factor (B)
- Data/parameter ratio >>1 is sought in X-ray crystallography
- Use the structure to infer biological/biochemical properties of the protein studied











Software for model building

- UCSF Chimera (rigid-body docking of structures into EM maps)
 - Visualizing "overall" cryoEM reconstructions
- Coot (rigid-body docking, manual model building, real-space refinement)
 - Visualizing "local" segment of cryoEM reconstructions
- Refmac (refinement of atomic coordinates)
- Phenix (model building "map_to_model" derived from resolve & refinement of atomic coordinates)
- Buccaneer (model building)
- PathWalker(model building)
- DireX (flexible fitting/refinement of atomic coordinates)
- MDFF (flexible fitting/refinement of atomic coordinates)
- Rosetta (model building & refinement of atomic coordinates)

• ..



Use resolution of high-resolution structures to extend structure determination at lower resolution



De novo structure determination

- Obtaining a structure with no other information than the sequence of the protein(s) and a cryoEM reconstruction
- Challenging at 3-5 Å resolution
- Hand tracing is time consuming
- Hand tracing is sometimes not possible
- Why do crystallographic model building softwares fail?
 - Usually work well at resolution ≤3 Å
 - Assign sequence primarily using side chain density

Take advantage of knowledge-based sequence/structure information:

Local sequence confers conformational preferences



For <u>9-residue windows</u> centered at each residue position, representative backbone conformations (<u>fragments</u>) are predicted using local sequence.

Wang RYR et al (2015) Nature Methods

Use fragment consistency to select correct fragments



$$score_{total}(F) = w_{dens} \sum_{f_i \in F} score_{dens}(f_i)$$
$$+ w_{overlap} \sum_{f_i, f_j \in F} score_{overlap}(f_i, f_j)$$
$$+ w_{close} \sum_{f_i, f_j \in F} score_{close}(f_i, f_j)$$
$$+ w_{clash} \sum_{f_i, f_j \in F} score_{clash}(f_i, f_j)$$

- Near-native fragments match the density well and are consistent with one another:
 - Non-clashing
 - Assign the same residue to the same position
 - Nearby residue in sequence are nearby in cartesian space
- Score function evaluates the consistency of a set of fragments
- Monte-Carlo sampling finds fragment set optimizing the score function.
- Works up to 4.8Å resolution

Wang RYR et al (2015) Nature Methods

Prefusion MHV spike glycoprotein 4.0Å resolution

Postfusion MHV spike glycoprotein 4.1Å resolution



Walls AC et al (2016) Nature Walls AC et al (2017) PNAS

Lexi Walls



 In placement: long fragments do not necessarily capture structure accurately in regions of high local variation



 In rebuilding: sampling long segments requires all the residues to be correct in order to see an energy signal



- Use shorter fragments for more accuracy in variable regions while still maintaining information on sequence preference
- Employ a greedy conformational strategy to handle the large conformational space
- Model completion

Works with unsegmented density & symmetry



Frank DiMaio

Brandon Frenz

Frenz B et al (2017) Nature Methods

- Reduction of search space leads to improvements
 - Penalizing density discontinuities



• Explicit modeling of β-sheets



Rotamer-like density matching



Frenz B et al (2017) Nature Methods



Bfactor

Frenz B et al (2017) Nature Methods



Frenz B et al (2017) Nature Methods Zhang X et al (eLIFE2013)

Prefusion
MHV spike
glycoprotein
4.0Å resolution



Glycosylation site (confirmed using MS)

Frenz B et al (2017) Nature Methods Walls AC et al (2016) Nature

Prefusion **HCoV-NL63** spike glycoprotein 3.4Å resolution Glycosylation sites (confirmed using MS)

Frenz B et al (2017) Nature Methods Walls AC et al (2016) Nature Struct Mol Biol

Rosetta glycan refinement





Frank DiMaio



Brandon Frenz

PDCoV

78 N-linked glycans

Alex Xiong

Lexi Walls



Automated glycan detection Refinement using stereochemical restraints (chair conformation)

102 N-linked glycans

Walls AC et al (2016) Nature Struct Mol Biol Xiong X et al (2017) J Virol Frenz B et al, unpublished



correspondence

Carbohydrate anomalies in the PDB

To the Editor: The importance of carbohydrates both to fundamental cellular biology and as integral parts of therapeutics (including antibodies) continues to grow. The presence of the correct glycans is important for the beneficial effects of therapeutic glycoproteins and is likely to be increasingly required by regulatory agencies. However, carbohydrates (and other small molecules) are handled poorly in macromolecular structural biology. When such small molecules are present in macromolecule structures, they are often reported with stereo- and regiochemical errors and in unlikely conformations. Stereoand regiochemistry should always be correct, and although conformational distortions may reflect interactions taking place in a complex¹, most are also likely to be erroneous-resulting from poor chemical understanding and lack of appropriate stereochemical restraints in refinement, often against low-resolution data².



Figure 1 | Distribution of D-pyranoside ring conformations as a function of resolution for all N-linked sugars (at distance <2.0 Å) in the PDB as of January 2015, identified by their Chemical Component Dictionary IDs: NAG, NDG, MAN, BMA, BGC, GLC, GAL and GLA. E/H, envelopes and half-chairs; B/S, boats and skew-boats; wavy lines denote the main ring plane. For clarity, an envelope is depicted at $\theta = 45^{\circ}$ and a half-chair at $\theta = 135^{\circ}$, and skew-boat is omitted from the equator.

Rosetta glycan refinement



Walls AC et al (2016) Nature Struct Mol Biol

Rosetta comparative modeling (RosettaCM)

- Structure(s) sharing sequence similarity to the target structure
 - Partial structures
 - Multiple structures
- Fragments predicted using local sequence
- Tuning of the frequency with which the two sources of information are used
- Also allows model completion when 70% built model available (although RosettaES outperforms RosettaCM in most cases)
- Recovers core packing from crystal structures

Song Y et al (2013) Structure



Rosetta comparative modeling (RosettaCM)





James ZM, Borst AJ et al (2017) PNAS

Rosetta comparative modeling (RosettaCM)



Deleted disordered loops between transmembrane helices & Truncated residue to Cβ in the C-terminal CNBD due to dampened resolution

James ZM, Borst AJ et al (2017) PNAS

Rosetta density-guided iterative refinement

- Fragments predicted using local sequence
 - Size can be tuned: the longer the fragment, the more divergence is allowed
 - Targets:
 - Regions with conformational strains and poor fit to density
 - Random segments of the chain(s)
 - User-defined regions
- 100-1000 models



- Produces atomic-level accuracy models using maps determined at 4.5 Å resolution or better
- Large radius of convergence
- Voxel size refinement
- B-factor refinement after model completion

DiMaio F et al (2015) Nature Methods Wang RYR et al (2015) Nature Methods

Rosetta density-guided iterative refinement



Computational cost of running Rosetta

- Several hundreds to thousands of trajectories
 - One cpu/trajectory
 - Selection based on Rosetta energy function and/or map/model FSC
- Run time depends of size of model to be built
 - 1.5h for a 30-residue fragment for RosettaES using 16 cores
 - 2-4h/model for RosettaCM, Rosetta density-guided iterative refinement (75-120 kDa/protomer applying symmetry)
 - 15 min/model for Rosetta relax
- Backfill queue
 - Makes use of free cpu cycles
 - Typical of super-computer centers

Rosetta structure determination pipeline

- Rosetta energy function and conformational sampling are valuable at bridging the resolution gap between for near-atomic resolution reconstructions
- Automatic de novo structure determination at 3.5-4.5Å is possible (up to 4.8A!)
- Versatile pipeline

Model building using several maps

- Different subsets of particles (classification)
- Different softwares



Cheng Y et al (2015) Nature Methods



Which maps should we deposit?

- Recent 3dem/ccpem discussion
- Full map sharpened (current standard)
- Unsharpened map (TRPV1, LliK...)
 - Caveat: end-user needs got know how to sharpen a map
- Half maps
- Mask used for resolution estimation

- The model should be stereochemically sound (Molprobity, Privateer)
 - http://molprobity.biochem.duke.edu/
 - Ramachandran statistics
 - Clashes (steric overlaps)
 - Distribution of favored rotamers

All-Atom Contacts	Clashscore, all atoms:	1.2		$100^{\text{th}} \text{ percentile}^* (N=744, 1.94 \text{\AA} \pm 0.25 \text{\AA})$		
	<u>Clashscore</u> is the number of serious steric overlaps (> 0.4 Å) per 1000 atoms.					
Protein Geometry	Poor rotamers	0	0.00%	Goal: <1%		
	Ramachandran outliers	0	0.00%	Goal: <0.05%		
	Ramachandran favored	151	97.42%	Goal: >98%		
	MolProbity score	0.95		$100^{\text{th}} \text{ percentile}^* (N=11856, 1.94 \text{\AA} \pm 0.25 \text{\AA})$		
	<u>Cβ</u> deviations >0.25Å	0	0.00%	Goal: 0		
	Bad backbone bonds:	0/1288	0.00%	Goal: 0%		
	Bad backbone angles:	0/1734	0.00%	Goal: <0.1%		

All-Atom	Clashscore, all atoms:	135.29		2 nd percentile [*] (N=37, 3Å - 9999Å)		
Contacts	Clashscore is the number of serious steric overlaps (> 0.4 Å) per 1000 atoms.					
	Poor rotamers	4681	46.69%	Goal: <0.3%		
	Favored rotamers	3431	34.22%	Goal: >98%		
	Ramachandran outliers	3282	28.73%	Goal: <0.05%		
Protein	Ramachandran favored	4982	43.62%	Goal: >98%		
Geometry	MolProbity score^	4.86		1^{st} percentile [*] (N=342, 3.50Å ± 0.25Å)		
	Cβ deviations >0.25Å	9	0.08%	Goal: 0		
	Bad bonds:	7 / 94333	0.01%	Goal: 0%		
	Bad angles:	490 / 128344	0.38%	Goal: <0.1%		
Peptide Omegas	Cis Prolines:	50 / 707	7.07%	Expected: ≤ 1 per chain, or $\leq 5\%$		

Chen VB et al (2010) Acta Cryst D

Barad BA et al (2015) Nature Methods

- EMRinger $(\chi 1)$
- http://www.ccp4.ac.uk/html/privateer.html



Agirre J et al (2015) Nature Struct Mol Biol

• The model should agree with the density



• Half maps are useful



Campbell MG et al (2015) eLIFE

- The model should be **<u>cross-validated</u>**
 - Using information that was not used during refinement!
- X-ray crystallography: reciprocal space data
 - Set aside a few percent of reflections (test set: Rfree)
 - Use the rest for model refinement (Rwork)
 - Measures agreement between Rwork and Rfree
- CryoEM data: real space data
 - Gold-standard refinement divides dataset in two halves
 - One half used as training (work) map
 - Other half used as testing map (free)
- Last iteration uses all the data (minimization only!)

Brunger AT (1992) Nature



DiMaio F et al (2013) Protein Science

- At low resolution (3-5Å), any additional source of information is most welcome!
- Glycosylation sites

Walls AC et al (2016) Nature Struct Mol Biol

Disulfide bonds

Walls AC et al (2017) PNAS

Cross-linking/MS (can be incorporated in RosettaCM)



- Coevolution encodes structural information
- Contacts in proteins are evolutionary conserved and encoded in a multiple sequence alignment due to coevolution
- By measuring coevolution, one can infer contacts in proteins



Ovchinnikov S et al (2014) eLIFE Ovchinnikov S et al (2017) Science

Sergey Ovchinnikov David Baker



Ovchinnikov S et al (2014) eLIFE Ovchinnikov S et al (2017) Science



- Ovchinnikov S et al (2017) Science
- Could provide an additional metric for cross-validation
 - Currently restricted to prokaryotic sequences

How to deal with uncertainty in model building

- Poorly ordered regions of a map can be:
 - Left unmodeled
 - Modeled and let B-factor account for it
 - Modeled but truncated at Cα or Cβ
 - Modeled with an occupancy of 0
 - A combination of the above
- Best strategy depends of the situation (*de novo* vs rebuilding)
- In any case, write in the manuscript what you have done

Model interpretation

- Model convergence
 - Rosetta ES/CM/density-guided iterative refinement
 - Mark Herzik & Gabe Lander's convergence server
 - https://doi.org/10.1101/128561
 - <u>http://www.lander-lab.com/convergence/</u>
- Agreement between local resolution and B-factor analysis



Walls AC et al (2016) Nature Struct Mol Biol

Acknowledgements

Veesler lab



Collaborators Frank DiMaio Brandon Frenz Felix Rey M. Alejandra Tortorici **Berend-Jan Bosch Bill Zagotta Zachary James David Baker** Sergey Ovchinnikov

NRAMM National Resource for Automated Molecular Microscopy



National Institute of General Medical Sciences



National Institute of Allergy and Infectious Diseases



NU

Rubicon







Ensemble Cap



Run time vs ensemble cap



Ensemble cap

Real-space *B*-factor refinement. To better model the density maps and generate more accurate models, we refined atomic *B* factors against the maps optimizing the real-space correlation between model and map. Given that atom *i* has a *B* factor B_i , we calculate the density of the model as

$$\rho_{\rm c} = \sum_{\rm atoms } i \left(\frac{\pi}{f_i + B_i / 4} \right)^{\frac{3}{2}} \exp \left(-\frac{\pi^2}{f_i + B_i / 4} \| x - x_i \|^2 \right)$$

Here, *f* is a scattering factor fit to each element. Our implementation makes use of a single-Gaussian scattering for each atom type, but it is straightforward to extend this to a standard five-Gaussian scattering model²⁷.

B-factor refinement is carried out using quasi-Newton optimization, with the gradient of the *B* factor of atom *i* (located at coordinates x_i) given in real space by

$$\frac{\partial \text{RSCC}}{\partial B_i} = \frac{1}{\sigma_c^2} \left(\sigma_c \frac{\partial \sum \rho_c \rho_o}{\partial B_i} - \sum \rho_c \rho_o \frac{\partial \sum \rho_c^2}{\partial B_i} \right)$$