### Post-Processing and Validation of 3D EM Reconstructions

- Model-based Refinement & the Need for Validation
- Validation at Low Resolution
- Tilt-Pair Validation
- Labels/Biochemistry
- Validation at High Resolution
- What should my map look like?
- Validating models built into EM maps
- Model stereochemistry

Peter Rosenthal

EMBO Image Processing 2019

Why is validation important?

- Map validation
- Model validation
- In course of structure determination, may indicate whether the experiment or data processing needs to be optimized.
- Data deposition
- Peer review
- Still a subject of research!

# Map Validation

- Is map correct (incorrect) at low resolution?
- What should the map look like at high resolution? (Contrast Restoration)
- Is resolution assessment exaggerated?
- Simple tests to demonstrate validity of map
- User still decides how to process the data



#### Outcome of the First Electron Microscopy Validation Task Force Meeting

Richard Henderson,<sup>1</sup> Andrej Sali,<sup>2</sup> Matthew L Baker,<sup>3</sup> Bridget Carragher,<sup>4</sup> Batsal Devkota,<sup>5</sup> Kenneth H. Downing,<sup>6</sup> Edward H. Egelman,<sup>7</sup> Zukang Feng,<sup>5</sup> Joachim Frank,<sup>8,9</sup> Nikolaus Grigorieff,<sup>10</sup> Wen Jiang,<sup>11</sup> Steven J. Ludtke,<sup>3</sup> Ohad Medalia,<sup>12,21</sup> Pawel A. Penczek,<sup>13</sup> Peter B. Rosenthal,<sup>14</sup> Michael G. Rossmann,<sup>15</sup> Michael F. Schmid,<sup>3</sup> Gunnar F. Schröder,<sup>16</sup> Alasdair C. Steven,<sup>17</sup> David L. Stokes,<sup>18</sup> John D. Westbrook,<sup>5</sup> Willy Wriggers,<sup>19</sup> Huanwang Yang,<sup>5</sup> Jasmine Young,<sup>5</sup> Helen M. Berman,<sup>5</sup> Wah Chiu,<sup>3</sup> Gerard J. Kleywegt,<sup>20</sup> and Catherine L. Lawson<sup>5,\*</sup>

Structure. 2012 Feb 8;20(2):205-14. doi: 10.1016/j.str.2011.12.014.

#### MODEL-BASED DETERMINATION OF ORIENTATION PARAMETERS

NEW MODEL



#### STARTING MODEL

(A) Four reference images (each 64  $\times$  64 pixels) used for picking from 1,024 random noise images (of 1,024  $\times$  1,024 pixels).



van Heel M PNAS 2013;110:E4175-E4177



model bias: the persistence of an incorrect map or map features during refinement

over-refinement or over-fitting: causes build-up of spurious noise features in the map during iterative refinement.

e.g. the map matches signal in the images at low resolution, but at high resolution predominantly matches noise.

over-refined maps can possess features that resemble side-chain densities but in incorrect locations.

### Cryo-EM image of a field of view of $\beta$ -galactosidase single particles (molecular weight, 450 kDa).



Henderson R PNAS 2013;110:18037-18041



### De novo determination of "Starting Map"

Map projections agree with individual raw images and class averages ("reference free") Distribution of particle orientations Symmetry Absolute Hand (Experimental Determination-Tilting)

#### **Other sources for starting map:**

Derived from another highly similar structure?

Density from X-ray model? Low-pass filtered

A spherical or cylindrical blob? Icosahedral or helical symmetry



Naydenova & Russo, Nature Communications, 2017

# Maximum Likelihood (ML) vs. CC (Align)

















Align 50

Structure

Average

First Ref.

Align 3





ML 274



Align 10



Structure

First Ref.

ML 10









а

#### Henderson et al. (2011) Wasilewski and Rosenthal (2014)



### **TILT AXIS FOR EACH PARTICLE PAIR**

### BEFORE OPTIMIZATION

AFTER OPTIMIZATION



#### Lu et al. Nature 2014

Tilt-pair for Hand Determination



### **Identifying Biochemical Labels in Maps**

#### Avidin/biotin



GFP

ATP synthase from Saccharomyces cerevisiae: location of subunit h in the peripheral stalk region.

Rubinstein et al. (2005) JMB

Structural basis for subunit assembly in the anaphase-promoting complex.

Schreiber et al. (2011) Nature

# MAP RESOLUTION





Resolution definition by separation of features. (a) When two points are far apart, there is a deep trough of density between them. (b) Two points are regarded as just resolved when the peak of one point spread function overlaps the first minimum of the other (Rayleigh criterion), see ref 60. (c) The point spread functions of two dots close together overlap to form one maximum, so that the points are not resolved.

The observed features of a map should be consistent with the resolution assessment









Map Resolution Should Be Reported, and Visible Structural Features Should Be in Accordance with the Claimed Resolution

Visibility of expected features – α-helices visible at 9 Å resolution? β-strands at 4.8 Å resolution? Side-chains beyond 4 Å?

Fourier Shell Correlation: Show the whole curve.



From: EMDataBank unified data resource for 3DEM Nucleic Acids Res. 2015;44(D1):D396-D403. doi:10.1093/nar/gkv1126 Nucleic Acids Res | © The Author(s) 2015. Published by Oxford University Press on behalf of Nucleic Acids Research.This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0/), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.



Methods Enzymol. 2010; 482: 73–100.

# What is a structure factor?





 $\sum F_1 \cdot F_2^*$  $F_1$ FSC  $\frac{1}{\sqrt{\sum |F_1|^2 \sum |F_2|^2}}$ N2  $F_2$  $\sum (\mathbf{S} + \mathbf{N1}) \cdot (\mathbf{S} + \mathbf{N2})^*$  $\sqrt{\sum |S + N1|^2 \sum |S + N2|^2}$  $\frac{\sum |S|^2}{\sum |S^2 + N^2|}$ FSC S~N FSC=0.5





#### Signal~Noise FSC=0.5

 $FSC_{full} = \frac{2FSC}{1+FSC}$ 

### Correlation with a perfect reference



Looks like figure-of-merit (Blow and Crick, 1959)

# Estimating C<sub>ref</sub>

$$C_{ref} = \sqrt{\frac{2FSC}{1 + FSC}}$$

C<sub>ref</sub>=0.5 FSC=0.14

### C<sub>ref</sub> corresponds to crystallographic FOM "Figure of Merit"

C<sub>ref</sub>=0.5 mean phase error 60 ° (last shell) interpretable by an atomic model

FSC	<b>FSC</b> <sub>FULL</sub>	C <sub>REF</sub>	<b>PHASE ERROR</b>	S/N <sub>1/2</sub>
0.50	0.67	0.82	35°	1.00
0.33	0.50	0.71	45°	0.71
0.14	0.25	0.50	60°	0.41





Resolution (1/Å)

# R<sub>work</sub>, R<sub>free</sub> in X-ray cross-validation to detect over-fitting



Free shells at high-resolution not included in refinement





#### Experimental results showing the utility of a tilted data collection strategy.



Dmitry Lyumkis J. Biol. Chem. 2019;294:5181-5197



#### HIGH RESOLUTION NOISE SUBSTITUTION

Chen et al. Ultramicroscopy 135 (2013) 24-35

Perform Single Particle EM analysis.



Repeat but substitute random phases beyond a selected resolution (HR-noise)

Any overfitted noise will show up as non-zero FSC.

Genuine information will show up as the area between the two curves

#### HIGH RESOLUTION NOISE SUBSTITUTION Chen et al. *Ultramicroscopy* **135** (2013) 24-35



# Assess overfitting

- Especially useful for new programs/methods
  - high-resolution noise substitution test



Chen et al. *Ultramicroscopy* **135** (2013) 24-35



Chen et al. *Ultramicroscopy* **135** (2013) 24-35



The difference between the FSC from the original data (FSCt) and the FSC with high-resolution noise substitution or phase randomization (FSCn) represents the true high-resolution signal in the map (FSCtrue)

Chen et al. Ultramicroscopy 135 (2013) 24-35

Resolution (1/Å)



Chen et al. Ultramicroscopy 135 (2013) 24-35

#### LOCAL RESOLUTION: RESMAP



#### RADIALLY-AVERAGED EM MAP AMPLITUDES COMPARED TO X-RAY



### **Contrast Restoration**

(Rosenthal and Henderson 2003)

Incorrect Scaling of High and Low Resolution Amplitudes makes map looks featureless

Problem: Application of a negative temperature factor amplifies both signal and noise

Suppress Noise by using Noise-Weighted Structure Factors C<sub>ref</sub>F Similar to Figure-of-Merit weighting









F



 $C_{ref}Fe^{(1000/4d2)}$ 



 $C_{ref}Fe^{(1000/4d2)}$ 

 $Fe^{(1000/4d2)}$ 







# Summary – Map Validation

- Recommend performing as many tests as possible to communicate the validity of a map. Furthermore, validation tools may be applied to optimizing map calculation or for testing new computational procedures.
- Heterogeneous particles raise several new problems in cryo-EM for which new validation tools may be required. High-resolution maps may be calculated from a small subset of the data (e.g. ~5%). Challenge!