3DEM Data Archiving & Validation, Challenges

Cathy Lawson

NYSBC CryoEM Course April 10, 2017



3DEM Data Archives for Structural Biology

PDB - coordinates - managed by wwPDB

EMDB – maps - managed by EMDataBank

EMPIAR – raw images - managed by EBI



Protein Data Bank (PDB)

Single global archive of 3-D macromolecular structures (>125,000 entries) est. 1971

1990s: First EM structures deposited

2017: ~1500 EM structures

Managed by the world-wide PDB organization



3

EM Data Bank (EMDB)

- Global archive for maps produced using 3DEM methods, est. 2002 @ EBI
- Since 2007, managed by EMDataBank partners, with support from NIGMS
- We expect to reach 5000 entries by end of 2017





EMPIAR (pronounced empire)

Archive for raw electron microscopy image data est. @ EBI/PDBe in 2013

- ~100 entries
- Supported by MRC, BBSRC



EMDataBank.org

Global portal for deposition and retrieval of 3DEM density maps, atomic models, and associated metadata

Resource for news, events, software tools, data standards, validation methods for the 3DEM community



News

Remediation of 3DEM Entries

EMDataBank/Unified Data Resource

in the Protein Data Bank

The wwPDB and the

All news

Unified Data Resource for 3-Dimensional Electron Microscopy

EMDataBank is a unified global portal for deposition and retrieval of 3DEM density maps, atomic models, and associated metadata, as well as a resource for news, events, software tools, data standards, validation methods for the 3DEM community.

For up-to-date information about map and model challenges, visit challenges.emdatabank.org.

Growth of 3DEM Archives





EM Structures @ 4 Å or better



EM Structures 2010 vs 2015

2010: Molecular Shapes

2015: Traceable Densities



0.5% of all entries in PDB (332 of 67500) 0.8% of all entries in PDB (905 of 112400)



3DEM Structure Validation



Importance of Validation

- J. Cohen, Is High-Tech View of HIV Too Good to Be True? Science 341, 443-444 (2013)
- R.M. Glaeser, Replication and validation of cryo-EM structures J. Struct. Biol. 184, 379-380 (2013)
- R. Henderson, Avoiding pitfalls of single particle cryo-electron microscopy: Einstein from noise, PNAS 110, 18037-41 (2013)
- M. van Heel, Finding trimeric HIV-1 envelope glycoproteins in random noise, *PNAS* 110, E4175-7 (2013)
- S. Subramaniam, Structure of trimeric HIV-1 envelope glycoproteins, PNAS 110, E4172-4 (2013)



EMD-5418 Y Mao, JG Sodroski *et al.* Molecular architecture of the uncleaved HIV-1 envelope glycoprotein trimer *PNAS* 110, 12438-12443 (2013)



Community Input for Validation

Task Force	Meeting/ Workshop	Chair(s)/Membership	Outcome
X-ray Validation Task Force	2008 (2015)	Randy Read (Univ of Cambridge) 17 members	(2011) <i>Structure</i> 19: 1395-1412
NMR Validation Task Force	2009, 2011, 2013 (x2), 2015	Gaetano Montelione (Rutgers) Michael Nilges (Institut Pasteur) 10 members	(2013) <i>Structure,</i> 21: 1563-1570
3DEM Validation Task Force	2010	Richard Henderson (MRC-LMB) Andrej Sali (UCSF) 21 members	(2012) <i>Structure</i> 20: 205-214
Small- Angle Scattering Task Force	2012, 2014	Jill Trewhella (Univ Sydney) 6 members	(2013) <i>Structure</i> 21: 875-881
Hybrid Methods Workshop	2014	Andrej Sali (UCSF), Torsten Schwede (Univ Basel), Jill Trewhella (Univ Sydney) 27 members	(2014) <i>Structure</i> 23: 1156-1167

12

EM Validation Task Force



Henderson *et al.* (2012) *Structure 20*, 205-214 <u>http://www.ncbi.nlm.nih.gov/pubmed/22325770</u>



EM VTF Recommendations

Main recommendations: EM maps

- Standards for assessing resolution and accuracy of a map need to be developed
- Structural features in a map should be in accordance with the claimed resolution
- Main recommendations: models fitted into EM maps
 - Criteria for assessing models in context of mapmodel fit need to be developed
 - Capability to archive coarse-grained representations of models is needed
- More research and development needed!



EM VTC on Resolution

- Resolution via FSC: "Deposition of a published map should include its full FSC curve to the Nyquist frequency on a linear spatial-frequency scale."
 - Masking: discouraged
 - Fully independent refinement: encouraged
- Low resolution structures (where secondary structure elements cannot be identified) can be validated using tilt-based methods



2015-Ongoing: Map, Model Challenges

- Goals: Develop benchmarks, encourage development of best practices in 3DEM reconstruction and model fitting, evolve criteria for validation, compare and contrast different approaches
- Developed by expert committees
- Results discussion via Participant Workshops/Journal Special Issues
- <u>http://challenges.emdatabank.org</u>



Benchmark Datasets

Map Challenge Targets: Raw Images @ EMPIAR



Model Challenge Targets: Maps @ EMDB





Challenges: Progress

Map Challenge:

27 participants submitted 66 maps Model Challenge:

16 participants submitted 106 models

Lots of Data Received! – assessments now in progress



3DEM Structure Deposition





- New Version for depositing structures from X-ray, NMR, and EM Launched January 2016
 Deposit map to EMDB with associated model to PDB in single session
- Validation report provided



3DEM Deposition: Method

Experimental method

X-Ray Diffraction

http://deposit.wwp db.org/deposition/

- Electron Microscopy
 - Helical
 - Single particle
 - Subtomogram averaging
 - Tomography
- Solution NMR
- Neutron Diffraction
- Electron Crystallography
- Solid-state NMR
- Fiber Diffraction



3DEM Deposition: IDs

http://deposit.wwpdb.org/deposition/

Are you depositing coordinates with this submission?
No, experimental data only
Yes
Has the associated map been deposited previously?
No

O Yes

Requested accession codes

PDB
EMDB
BMRB



File uploads: 3DEM map/model submission in OneDep

✓Select file type...

0) Coordinates

Coordinates (mmCIF format) Coordinates (PDB format)

1) Main map (mandatory) EM map (MRC/CCP4 format)

2) Image for EMDB (mandatory)

Entry image for public display

3) Additional maps

Additional EM map (MRC/CCP4 format)

4) Masks

EM mask (MRC/CCP4 format)

5) Half (even-odd) maps

EM half map (MRC/CCP4 format)

6) Structure Factors

mmCIF (structure factors)

MTZ

Other Files

FSC file (XML format)

Ligand Image

FSC Curve Upload

Create xml format file using a software package (e.g., Relion, EMAN)
 Use PDBe's Server: PDBe.org/FSC





EM Validation Report

Value

30000

FSC 0.143

Not provided

JEOL 3200FSC

SINGLE PARTICLE

Source

Depositor

Depositor

Depositor

Depositor

Depositor

Depositor

19

Unified Data Resource fo

Property

Microscope

"Table 1"

Metrics of the

Reconstruction method

CTF correction method

Resolution determination method

Imposed symmetry

Number of images

Voltage (kV) 300 Depositor EM model Electron dose $(e^{-}/2)$ Not provided Depositor Minimum defocus (nm) Depositor 500Maximum defocus (nm) 2000 Depositor Magnification 50000 Depositor Image detector DIRECT ELECTRON DE-12 (4k x 3k) Depositor Percentile Ranks Metric Value Clashscore Ramachandran outliers 2.9% Sidechain outliers 0.1% Worse Better Percentile relative to all structures Percentile relative to all EM structures **EMDataBank**

EM Validation Report: Proposed additions

- Table of Map Parameters
- Map images
- Plots
- Model to map fit
- Example:

http://www.ebi.ac.uk/pdbe/entry/emdb/EMD -2984/analysis



Project Team



Baylor College of Medicine

Wah Chiu, PI Steven Ludtke Corey Hryc Grigore Pintilie Matthew Baker Matthew Dougherty **Rutgers University**

PROTEIN DATA BANK

Helen Berman, co-PI Catherine Lawson Raul Sala Brian Hudson John Westbrook

European Bioinformatics Institute

Protein Data Bank in Europe

ЈВе

Gerard Kleywegt, co-Pl Ardan Patwardhan Sanja Abbott Eduardo Sanz Garcia Ingvar Lagerstedt

Supported by NIH National Institute of General Medical Sciences

