# EMDataResource: Structure Data Archiving, Validation Challenges

Cathy Lawson EMDataResource & RCSB Rutgers University

NYSBC NCCAT Single Particle Short Course, March 5, 2020



### **Unified Data Resource for 3DEM**



Stanford University/SLAC



**Rutgers University** 



European Bioinformatics Institute

Established 2007 under NIH Support (R01GM079429) to:

- Develop Data Archives for 3DEM (EMDB + PDB)
- Promote Community Development of Validation and Standards



### **Project Website**

2019-07-10 : 8627 EMDB map entries, 3596 PDB coordinate entries RCSB PDB | PDBe



#### EMDataResource

Unified Data Resource for 3DEM

Home About - Deposit Search - Tools - Events - News Links Help -

Unified Data Resource for 3DEM

Global resource for 3-Dimensional Electron Microscopy (3DEM) structure data archiving and retrieval, news, events, software tools, data standards, validation methods, and community challenges

Quick link: model metrics challenge

#### **Recently released EMDB entries**

#### All recent entries



#### singleParticle 3.37 Å PCE-6oby Conformational states of Cas9sgRNA-DNA ternary complex in the presence of magnesium Zhu X, Clarke R, Puppala AK, Chittori S, Merk A, Merrill BJ, Simonovic M, Subramaniam S

EMD-0584

Deposited 2010-02-17

News

Search EMDB

All news

#### 2019 Model Metrics Challenge

EMDataResource is trying something new and has put together a smaller, shorter-timeline EM model challenge. The deadline for submission of completed models fitted to target maps (2-3 Å resolution range) is May-26 May 28 (3 PM US Eastern). Read more...

#### Reminder: EMDB Accession Codes are Changing Soon!

The allocation of 4 digits for EMDB accession codes will soon come to an end. New EMDB accession codes will include an additional digit and will



emdataresource.org

### **Growth of EM Archives**



updated weekly: emdataresource.org/statistics.html



### **EMDB** maps by year and resolution



updated weekly: emdataresource.org/statistics.html



### **Finding Cryo-EM Structures**

#### emdataresource.org/search.html

#### EMDB maps

#### **PDB EM models**

- EMDR Search
- Advanced Search @
   EBI
- RCSB-PDB
- PDBe

EMDB + PDB

- EM Navigator @ PDBj
- 3DBioNotes @ CSIC

#### Raw image data

• EMPIAR @ EBI



### **EMDR Search: Features**

Autocomplete: ribo

ribosome

#### Interactive keyword help:

?

deposited:[2002 TO 2004]





Sort,	configure,	filter,	download	l search	results:
,					

	fide Sea	rch Options			
Filter	8	Released Entries      Any Method     Any Sample Type	+ Any Resolu	tion Stat # U	pdate Results
Results       Entry Title       Entry Authors       Deposit Date       Header Release Date       Map Release Date       Processing Site       Method         Specimen type       Sample Molecular Weight       Microscope Model       Detector       Energy Filter       Electron Dose         Number of Particles       Imposed Symmetry       Reported Resolution       Fitted Model PDB       Citation Authors       Citation Title         Citation Year       Journal Name       PubMed       DOI       Funding					
Found	1113 EM	IDB map entries for query ribosome (Released Entries, Any Method, Any Sample Type, Any Specir	men Type, Any Re	isolution Status)	1
Copy	CSV	Excel Column visibility * Show 10 rows *	Pinds No	ouno.	
* 1D-	Status	Entry Title	Deposit ** Date	Resolution == (Å)	Associated % PD6
* 4D- 147	Status REL	1 Entry Title Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP(sordarin)	Deposit Date	Resolution % (Å) 4.4	Associated 1% PDB 6gq1
4D- 047 048	Status REL REL	Entry Title Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP/sordarin) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GDP+AIF4/sordarin)	Deposit 10 Date 2018-06-07 2018-06-07	Resolution % (Å) 4.4 3.9	Associated % PDB 6gq1 6gqb
40- 047 048 049	REL REL REL	Entry Title Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP)sordarin) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GDP+AJF4/sordarin) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP)	Deposit Date 2018-06-07 2018-06-07 2018-06-08	Resolution % (Å) 4.4 3.9 4	Associated % PDB 6gq1 6gqb 6gqy
MD- 047 048 049 055	Status REL REL REL REL	Entry Title Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP)sordarin) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GDP+AIF4/sordarin) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GMPPCP) Cryo-EM reconstruction of yeast 80S ribosome in complex with mRNA tRNA and eEF2 (GDP+AIF4/sordarin): SSU-focused map	Deposit Date 2018-06-07 2018-06-07 2018-06-08 2018-06-14	Resolution % (Å) 4.4 3.9 4 4.3	Associated % PDB 6gq1 6gqb 6gqv

#### **EMDR Search: Demo**

#### emdataresource.org/solrsearch.html

#### What EMDB map entries would you like to find?



#### **EMDR Search: Programmatic Access**

<u>https://www.emdataresource.org/node/solr/emd/</u> <u>select?q=5127&fl=id,PDB&rows=10</u> retrieves:

{"responseHeader":{ "status":0, "QTime":3, "params":{ "q":"5127", "fl":"id,PDB", "rows":"10"}}, "response":{"numFound":1,"start":0,"docs":[ { "PDB":["3iyd"], "id":"5127"}] }}



### **Cryo-EM Structure Deposition**

EMDB, PDB



# wwPDB *meDep* System

- Deposition system for X-ray, NMR, and EM Structures
- EM Depositions:
  - Map to EMDB
  - Coordinate model to PDB
- Validation report is produced



# File uploads for EM deposition

✓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 XML format file
- Create XML via a software package (e.g., Relion, EMAN, cisTEM), or use PDBe's <u>FSC server</u>



## **PDB Deposition Policies**

- Map deposition to EMDB is mandatory for PDB depositions of 3DEM atomic coordinates
- MX atomic coordinates must be deposited in PDBx/mmCIF format
  - PDB format is still accepted for 3DEM and NMR (will change in future)
- Coordinates and experimental data share same release status (REL, HPUB, or HOLD).
  - Coordinate and experimental data are released simultaneously.
- ID's issued only after mandatory metadata are provided
- PI contact information and ORCiD ID are mandatory
- Author-initiated coordinate versioning is allowed post release



wwpdb.org/documentation/policy

# **3DEM Metadata Collection**

- High-level classification of the EM experiment
- Software used for data collection, data processing, data analysis, structure calculations, and refinement
- Sample description (e.g., assembly, virus)
- Data collection (e.g., diffraction, imaging)
- Sample preparation (e.g., specimen, buffer, tomography, condition)
- Image processing and reconstruction
- Structure analysis for 3D fitting of atomic coordinates



### **mmCIF Data Dictionary for 3DEM**

Top Level	Sample/Specimen	Image Processing &	Experimental Data
em_experiment	Preparation	Reconstruction	em_map*
em_software	em_buffer_	em_3d_reconstruction_	em_structure_factors*
	em_buffer_component	em_image_processing	<u>em_layer_lines</u> *
Sample Description	em_crystal_formation	em_particle_selection	
em_entity_assembly_	em_embedding	em_volume_selection	
em entity assembly molwt	em_sample_support_	em_ctf_correction	
em_entity_assembly_naturalsource	em_specimen		
em_entity_assembly_recombinant	em_staining	em_2d_crystal_entity	
<u>em_virus_entity_</u>	em_vitrification_	em_3d_crystal_entity	
<u>em_virus_natural_host</u>		em_helical_entity_	
<u>em_virus_shell</u>	em_fiducial_markers*	em_single_particle_entity_	
	em_focused_ion_beam*		
Data Collection	<u>em_grid_pretreatment</u> *	em_euler_angle_assignment*	
em diffraction	em_high_pressure_freezing*	em_final_classification*	
em_diffraction_shell	em_shadowing*	em_start_model*	
em_diffraction_stats	em_support_film*		
em_image_recording	em_tomography*	Structure Analysis	
em_image_scans_	em_tomography_specimen*	em_3d_fitting	
em_imaging_	em_ultramicrotomy*	<u>em_3d_fitting_list_</u>	
em_imaging_optics		em_fsc_curve*	

\*All categories are collected by the OneDep system. Data from most categories are archived in both PDB and EMDB; asterisked categories are archived only in EMDB.



#### **Cryo-EM Structure Validation**



# Validation Report for EM Structures v1 (2016-2019)

- Map resolution reported by depositor
- Model geometry statistics
- No fit-to-map validation





# Validation Report for EM Structures v2 (2020-)

#### New:

- Map, Map+Model Images
- FSC curve(s)
- Rotationally averaged power spectrum
- Fit-to-Map: Atom inclusion at recommended contour level



#### EMD-0273

#### Map Images

6.1 Orthogonal projections (i)



6.4 Orthogonal surface views (i)



### **Resolution Estimate by FSC**



Blue: FSC curve; Vertical Black Line: reported resolution



EMD-0273

### **Atom Inclusion**





EMD-0273 PDB 6UH7





### **Validation Challenges**



# EM Validation Task Force 2010 Recommendations

- Full FSC curve from independent half-maps
- Model Stereochemistry same as X-ray / NMR
- Other Metrics: More
   Research Needed



Henderson *et al.* (2012) *Structure 20*, 205-214 <u>10.1016/j.str.2011.12.014</u>



# Validation in a Changing Landscape



- How accurate are the maps?
- Do the models conform to good valence geometry and stereochemistry?
- How well do the models fit the maps?
- What are the metrics for evaluation and are they good enough?

To answer these questions, we conducted a series of Challenges



#### **Promoting Standards Development: Challenges**

Challenge Workshop(s)	Resol (Å)	Goals for Participants and Assessors		
2010 Model Challenge 2011 Hawaii 2012 Houston	2.5-24	<ul> <li>Produce best models against selected maps</li> <li>Explore segmentation, secondary structure detection, rigid body, flexible fitting, <i>ab initio</i></li> </ul>		
2016-2017 Map Challenge 2017 Lake Tahoe 2017 Stanford/SLAC		<ul> <li>Produce best maps from selected raw images</li> <li>Produce best models against selected maps</li> </ul>		
2016-2017 Model Challenge 2015 Boston 2017 New Orleans 2017 Stanford/SLAC	2.5-5	<ul> <li>Compare reconstruction, modeling practices</li> <li>Explore assessment strategies esp. map and model fit-to-map</li> </ul>		
2019 Model "Metrics" Challenge 2019 Stanford/SLAC	1.8-3.1	<ul> <li>Produce best models against selected maps</li> <li>Explore Model metrics with focus on Fit-to-Map</li> </ul>		

Unified Data Resource for 3DEM

# **2010 Model Challenge: Participants**

Chandrajit Bajaj, David Baker, Mariah Baker, Matthew Baker, Helen Berman, Radhakrishna Bettadapura, Virginia Burger, Kwok-Yan Chan, Chakra Chennubhotla, Wah Chiu, Frank DiMaio, Joachim Frank, Yaser Hashem, Tommy Hofmann, Shuiwang Ji, Tao Ju, Mert Karakaş, Steffen Lindert, Steven Ludtke, Cathy Lawson, Gerard Kleywegt, Jing He, Corey F. Hryc, Jens Meiler, Kamal Al Nasr, Grigore Pintilie, Ian Rees, Eduard Schreiner, Gunnar F. Schröder, Klaus Schulten, Dong Si, Phoebe L. Stewart, Leonardo G. Trabuco, Zhe Wang, Nils Wötzel, Qin Zhang





# **2010 Model Challenge: Observations**

- Established community around a common problem
- Identified critical standardization issues related to data deposition
- Identified issues to explore in future challenges
- Identified Challenges as mechanism to establish modeling benchmarks





#### 10.1002/bip.22081



#### 2016-2017 Map & Model Challenges: Participants

Maps: Committee: B Carragher (Chair), J-M Carazo, W Jiang, J Rubinstein, P Rosenthal, F Sun, J Vonck Data Contributors: Y Deng, F Sun, M Campbell, B Carragher, C Russo, L Passmore, J-P Armache, M Liao, Y Cheng, X Bai, S Scheres, Z Wang, W Chiu, A Bartesaghi, S Subramaniam Challenge Participants: A Punjani, A Nans, A Leith, A Chakraborty, JB Heymann, CO Sorzano, C Gati, D Tegunov, D-H Chen, F Li, G Yu, JM Bell, J Chen, JG Montoya, J Gomez-Blanco, K Yang, L Donati, L Estrozi, M Nilchian, N Caputo, N Grigorieff, S Stagg, S Shakeel, S Scheres, S Ludtke, X Bai, Y Lu Assessors: M Holmdahl, A Patwardhan, R Marabini, J-M Carazo, JB Heymann, J Mendez, S Stagg, G Pintilie,W Chiu, S Jonic, E Palovcak, J-P Armache, J Zhao, Y Cheng

**Models**: **Committee**: P Adams (Chair), A Brunger, R Read, T Schwede, M Topf, G Kleywegt **Data Contributors**: S Fromm, C Sachse, J-P Armache, Y Cheng, M Campbell, B Carragher, S-H Roh, C Hryc, W Chiu, Z Wang, A Bartesaghi, S Subramaniam, Xi Bai, S Scheres, N Fischer, H Stark, W Li, Z Liu, J Frank **Challenge Participants**: A Joseph, M Topf, B Frenz, F DiMaio, B Mao, C Hryc, W Chiu, D Kihara, G Terashi, D Matthies, G Schroeder, T Braun, K Wang, I Yu, H Zhou, M Baker, P Afonine, R McGreevy, A Singharoy, S Chittori, S Grudinin, A Hoffmann, T Kawabata, H Nakamura, T Terwilliger, T Croll, W Cao **Assessors**: A Krystafovych, A Joseph, M Topf, J Richardson, T Terwilliger, B Barad, J Fraser, A Jakobi, C Sachse



#### 2016-2017 Map & Model Challenges: Observations

- Innovative methods for map and model fit-to-map assessment were introduced
- Map quality depended on participant level of experience
- Maps reported as the same resolution looked different from each other
- Models were "all over the place"



10.1016/j.jsb.2018.10.004



#### 2016 Map & Model Challenges: Recommendations

- Map Resolution by independent half-map FSC: uniform definition + software implementation needed
- Novel model-based methods may be useful for estimating map resolvability
- Needed: further review of fit-to-map metrics



#### 2019 Model "Metrics" Challenge: Participants

Paul Adams, Pavel Afonine, Matt Baker, Helen Berman, Paul Bond, Tom Burnley, Renzhi Cao, Jianlin Cheng, Wah Chiu, Grzegorz Chojnowski, Kevin Cowtan, Frank DiMaio, Dan Farrell, James Fraser, Mark Herzik, Soon Wen Hoh, Maxim Igaev, Agnel Joseph, Daisuke Kihara, Andriy Kryshtafovych, Dilip Kumar, Cathy Lawson, Shanshan Li, Sumit Mittal, Bohdan Monastyrskyy, James Murray, Mateusz Olek, Colin Palmer, Ardan Patwardhan, Greg Pintilie, Alberto Perez, Jane Richardson, Peter Rosenthal, Daipayan Sarkar, Luisa Schaefer, Mike Schmid, Gunnar Schröder, Mrinal Shekhar, Dong Si, Abishek Singharoy, Genki Terashi, Tom Terwilliger, Maya Topf, Andrea Vaiana, Liguo Wang, Christopher Williams, Martyn Winn, Xiaodi Yu, **Kaiming Zhang** 







Unified Data Resource for 3DEM

### **2019 Challenge Targets**





## **Challenge** Pipeline



model-compare.emdataresource.org



Cc	Backbone	Conformation	Backbone	CaBLAM Co-trace Co-only virtual dihedrals CaBLAM Conformation Co and CO-containing virtual dihedrals MOLPROBITY Ramachandran	
rdinates	Sidechain		Sidechain	MOLPROBITY Rotamer	
Only V	Valence Geometry		ry	PHENIX Bond   Bond angle   Chirality   Planarity   Dihedral	
C	Clashes			MOLPROBITY Clashscore	
E		Energy		PROQ3 energy and predicted features	
	Full map density		Full map density	TEMPY CCC   PHENIX box_CC	
Cc	Density within a m	Correlation	Density within a mask	TEMPy CCC_overlap   Segment Mander's Overlap PHENIX CC_peaks   CC_volume   CC_mask	
	Density-derived fu		Density-derived functions	TEMPY Mutual Information(MI)   MI_overlap   Laplacian Filtered	
to Map	Density at atom po		Density at atom positions	MAPQ Q-score: vs Reference Gaussians (r=0-2 Å)	
E	Single point	FSC curve	Single point	PHENIX Resolution Map-Model FSC = 0.5	
F	Integration		Integration	CCPEM REFMAC5 FSCavg curve area to defined resolution limit	
A	Atom Inclusion			TEMPy Envelope   EMDB Atom Inclusion	
R	Rotamer			EMRinger Z-score protein $C_{\gamma}$ atom paths around $\chi 1$	
	Ca Superposition		Ca Superposition	OPENSTRUCT RMSD-Ca	
SI	Distance cutoffs	Superposition	Distance cutoffs	OPENSTRUCT Global Distance Calculation (GDC) all   sidechain Global Distance Test (GDT) total score   high accuracy	
eterence Iodel	Sequence assignr		Sequence assignment	PHENIX seq match   Co atom position match   overall score	
	*Multiple reference		*Multiple references	DAVIS-QA average of pairwise GDT_TS scores	
	Per chain	Distances	Per chain	LDDT Local difference distance test	
	All chains		All chains	OPENSTRUCT oligomeric LDDT   weighted oligomeric LDDT	
Models	Contact area	Contacts	Contact area	CAD Contact Area Difference	
sensus*	Shared contacts		Shared contacts	OPENSTRUCT Quaternary Structure (QS) best, global	
	Hydrogen bonds		Hydrogen bonds	HBPLUS H-bond Precision all   nonlocal   Similarity all   nonlocal	
to Map	Density-derived fu Density at atom po Single point Integration Ca Superposition Distance cutoffs Sequence assignt *Multiple reference Per chain All chains Contact area Shared contacts Hydrogen bonds	FSC curve Atom Inclusion Rotamer Superposition Distances Contacts	Density-derived functions Density at atom positions Single point Integration Ca Superposition Distance cutoffs Sequence assignment *Multiple references Per chain All chains Contact area Shared contacts Hydrogen bonds	TEMPY Mutual Information(MI)   MI_overlap   Laplacian Filtered MAPQ Q-score: vs Reference Gaussians (r=0-2 Å) PHENIX Resolution Map-Model FSC = 0.5 CCPEM REFMAC5 FSCavg curve area to defined resolution lin TEMPy Envelope   EMDB Atom Inclusion EMRinger Z-score protein $C_{\gamma}$ atom paths around $\chi 1$ OPENSTRUCT RMSD-Ca OPENSTRUCT Global Distance Calculation (GDC) all   sidech Global Distance Test (GDT) total score   high accuracy PHENIX seq match   Ca atom position match   overall score DAVIS-QA average of pairwise GDT_TS scores LDDT Local difference distance test OPENSTRUCT oligomeric LDDT   weighted oligomeric LDDT CAD Contact Area Difference OPENSTRUCT Quaternary Structure (QS) best, global HBPLUS H-bond Precision all   nonlocal   Similarity all   nonloc	

## 2019 Challenge: 4 Metrics Stood Out

Coordinates alone: **CaBLAM**: virtual dihedrals based on Ca, C=O compared to statistics from high quality PDB models [Williams 2018]

**MAPQ Q-score**: Per-atom Correlation vs Reference Gaussian (r=0-2 Å) [Pintilie in press]

**PHENIX Resolution Map-Model FSC** = 0.5 [Afonine 2018]

**EMRinger Z-score** "rotamericity" of map for protein  $C_{\gamma}$  atom paths around  $\chi 1$  [Barad 2015]

Fit-to-Map:

### **2019 Challenge: Observations**

- ab initio methods represented in the challenge performed extremely well for near-atomic maps (1.8 - 3.1 Å)
- For evaluating conformation: CaBLAM was a useful "orthogonal" metric to Ramachandran statistics
- Within single map targets, all fit-to-map metrics were equivalent (similar model rankings)
- only Q-score, EMRinger, and Map-Model FSC @ 0.5 provided useful comparisons across map targets



### **2019 Challenge: Recommendations**

- Most fit-to-map metrics are fine for optimization against a single experimental map
- Resolution-sensitive metrics are preferred for ranking diverse structures in an archive
- CaBLAM is a valuable new tool for evaluating protein backbone conformation



### **Topics for Future Challenges**

- Lower Resolution
- Membrane Proteins
- Nucleic Acids
- Models derived from Tomograms



### **Unified Data Resource for 3DEM**



# Stanford University (Baylor Coll. Med.)

Wah Chiu Greg Pintilie Mike Schmid

Steven Ludtke Matt Baker Corey Hryc Ian Rees

**UC Davis:** Andriy Kryshtafovych



#### **Rutgers University**

Cathy Lawson Helen Berman Brinda Vallat Brian Hudson

John Westbrook Batsal Devkota Raul Sala Chunxiao Bi



#### European Bioinformatics Institute

Ardan Patwardhan Gerard Kelywegt Sanja Abbott Ryan Pye Osman Salih Zhe Wang

Kim Henrick Richard Newman Christoph Best Glen van Ginkel Eduardo Sanz-Garcia Ingvar Lagerstedt





EMDataResource is funded by the US National Institutes of Health/National Institute of General Medical Science, **R01GM079429-12** 

### **EM Structure Validation Servers**

Мар:	Service/Name	Link
Overall Shape & Hand	Tilt-Pair	pdbe.org/tiltpair
Resolution by FSC	FSC	pdbe.org/FSC
Local Resolution	3DFSC	<u>3dfsc.salk.edu</u>
Local Resolution	Scipion	scipion.cnb.csic.es/m/myresmap#

Model:	Service/Name	Link
Stereochemistry, compare with all PDB structures	wwPDB	validate.wwpdb.org
Stereochemistry	Molprobity	molprobity.biochem.duke.edu
Nucleic Acid conformation	DNATCO	dnatco.org

Map/Model Fit:	Service/Name	Link
"backbone bumpiness"	EMRinger	emringer.com (@UCSF)

See also: <u>www.emdataresource.org/validation.html</u>



### References

- Lawson CL, Chiu W (2018) Comparing cryo-EM structures (Editorial). J Struct Biol. 204, 523-526. <u>10.1016/j.jsb.2018.10.004</u>
- Patwardhan A & Lawson CL (2016). Databases and Archiving for CryoEM. Methods Enzymol 579, 393-412. <u>10.1016/bs.mie.2016.04.015</u>
- Lagerstedt I, et al (2013). Web-based visualisation and analysis of 3D electron-microscopy data from EMDB and PDB. J Struct Biol. 184, 173-81. <u>10.1016/j.jsb.2013.09.021</u>
- Henderson R, et al (2012) Outcome of the first electron microscopy validation task force meeting. Structure 20, 205-214. <u>10.1016/j.str.2011.12.014</u>

