# Protein Data Bank:

An open access resource enabling basic and applied research and education in biology and medicine

## John Westbrook, Ph.D.

RCSB PDB Data & Software Architect Lead





## **Overview**

- A bit of background about the PDB
- PDB data content
- PDB data representation and data quality standards
- The PDB biocuration platform
- Some key features delivered by the RCSB PDB



## Protein Data Bank

- First open access digital resource in biology (est. 1971 with 7 entries)
- Single global archive of 3-D macromolecular structures (contains >120,000 entries)
- Freely available to all at pdb.org
- US PDB headquartered at Rutgers/UCSD (NSF, NIH, DOE)
- US PDB part of Worldwide PDB with partners in EU and Japan





# Worldwide Protein Data Bank

- Established in 2003, wwPDB ensures data are freely & globally available in a common repository
- Collaborate on data quality and representation standards, and tools and procedures for biocuration
- Each partner delivers different services and views of the common repository of data



wwPDB Advisory Committee Meeting, 2015







## Technology



### **PDB** archive established

## Community



### CRYSTALLOGRAPHY Protein Data Bank

ta will be open tly by the Crystallograph Data Centre, Camb Brookhaven National Laboratory The system will be responsible for oring atomic coordinates, structure





**IUCr** deposition guidelines





Standardization,

RCSB

robots



**Experimental data** required, wwPDB







### Validation standards

Metric	Percentile Ranks	Value
Rfree		0.189
tiscore		4
utliers		0
utiers		2,4%
utiers		0
Provide December	esistive to all X-ray structures relative to X-ray structures of similar resolution	Better
"Trany"		
The states	21	
10 10 1 10 10 10 10 10 10 10 10 10 10 10		

## **Diverse Molecular Content of the PDB**





## Supporting Rapidly Evolving Experimental Technologies



## X-ray Crystallography





Nuclear Magnetic Resonance Spectroscopy





R



# PDB Growth and Data Usage



As of 2015, ~55% increase in the number of global depositions since 2008

## PDB Depositors >800 new entries/month



PDB Users FTP and RSYNC Download Traffic in 2015: 526 million downloads



347 million





# PDB Data Content

- Atomic coordinates and primary experimental data
- Sample composition and preparation details
- Protein and nucleic acid polymer sequences and taxonomy
- Small molecules (ligands)
- Experimental data collection, structure solution, and structure refinement
- Structure classification by sequence, function and other criteria
- Citation and references to related data resources
- PDB entries contain 250 1200+ unique items of data

The 3.8 angstrom resolution cryo-EM structure of Zika virus. Sirohi, D., Chen, Z., Sun, L., Klose, T., Pierson, T.C., Rossmann, M.G., Kuhn, R.J. (2016) Science 352: 467-470 data\_5IRE

entry.id 5IRE

Shorf in Sing			5	LEA		2 m
em_entity_assembly.id	1		25	2.8-0		31:09
em_entity_assembly.details			A	and a	21.3	
Low passage Vero cells, derived f	rom Africa	n gr			To a series	A A A
were chosen for propagating the v	irus.		E			E.C.
em_entity_assembly.name	'Zika viru	s'	S.S.	1657	story.	Stor.
em_entity_assembly.type	VIRUS		CCC -		-Store	1.84
em_entity_assembly.assembly_id	1		548	CAS.	Contraction of	0
						1 mars
em_imaging.id		1	2966	a la	1. 6.10	Sealer -
em_imaging.entry_id		5IRE			1.00	1 Con
em_imaging.accelerating_voltage		300	8	-ain	and a	MAS
em_imaging.electron_source		'FIE		C.C.	Store 1	and the
em_imaging.illumination_mode		'FLO		- Mar	P.RU.	AND
em_imaging.microscope_model		'FEI		-	Mar Sec	
em_imaging.mode		BRI			and and	
em_imaging.nominal_cs		2.7				
em_imaging.nominal_defocus_max		2500				
em_imaging.nominal_defocus_min		1000				
em_imaging.nominal_magnification		14000	)			
em_imaging.recording_temperature_	maximum	80				
em_imaging.recording_temperature_	minimum	80				
em_imaging.specimen_holder_model		'FEI	TITAN	KRIOS	AUTOGRI	D HOLDER
em imaging.sample support id		1				



;

# loop\_

loop\_ \_entity\_poly.entity\_id \_entity\_poly.type \_entity\_poly.nstd\_linkage \_entity\_poly.nstd\_monomer \_entity\_poly.pdbx\_seq\_one\_letter\_code \_entity\_poly.pdbx\_seq\_one\_letter\_code\_can \_entity\_poly.pdbx\_strand\_id \_entity\_poly.pdbx\_target\_identifier 1 'polypeptide(L)' no no

; IRCIGVSRDFVEGMSGGTWVDVVLEHGGCVTVMAQDKPTVDIELVTTVSNMAEVRSYCYEASISDMASDSRCPTQGEA YLDKQSDTQYVCKRTLVDRGWGNGCGLFGKGSLVTCAKFACSKKMTGKSIQPENLEYRIMLSVHGSQHSGMIVNDTGHET DENRAKVEITPNSPRAEATLGGFGSLGLDCEPRTGLDFSDLYYLTMNNKHWLVHKEWFHDIPLPWHAGADTGTPHWNNKE ALVEFKDAHAKRQTVVVLGSQEGAVHTALAGALEAEMDGAKGRLSSGHLKCRLKMDKLRLKGVSYSLCTAAFTFTKIPAE TLHGTVTVEVQYAGTDGPCKVPAQMAVDMQTLTPVGRLITANPVITESTENSKMMLELDPPFGDSYIVIGVGEKKITHHW HRSGSTIGKAFEATVRGAKRMAVLGDTAWDFGSVGGALNSLGKGIHQIFGAAFKSLFGGMSWFSQILIGTLLMWLGLNTK NGSISLMCLALGGVLIFLSTAVSA

; IRCIGVSNRDFVEGMSGGTWVDVVLEHGGCVTVMAQDKPTVDIELVTTTVSNMAEVRSYCYEASISDMASDSRCPTQGEA YLDKQSDTQYVCKRTLVDRGWGNGCGLFGKGSLVTCAKFACSKKMTGKSIQPENLEYRIMLSVHGSQHSGMIVNDTGHET DENRAKVEITPNSPRAEATLGGFGSLGLDCEPRTGLDFSDLYYLTMNNKHWLVHKEWFHDIPLPWHAGADTGTPHWNNKE



# **Community Data Standards**

- PDB manages data using the macromolecular extension to the Crystallographic Information Framework (mmCIF) originally developed as an IUCr data standard
- PDB coordinates the extension of the standard (PDBx/mmCIF) to support the broader needs of both contributors and users of the archive (> 4400 data terms)
- A PDBx/mmCIF Working Group of community experts and methods developers oversees the evolution of the standard and ensures that the standard is well supported by key community software tools.
- PDB hosts community workshops to support the data standard and maintains a web site serving PDBx/mmCIF data dictionaries, schema and software tools (mmcif.wwpdb.org)





# PDBx/mmCIF Development Timeline



DATA BANK

# **Community Standards for Data Quality**

**Method-specific Community Validation Task Forces** have been convened to collect recommendations and develop consensus on data quality standards, identify software tools to perform required validation tasks, and to define related content requirements for archiving.



Task Force	Meeting/ Workshop	Chair(s)/Membership	Outcomes	
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 2016	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	2011 2014	Jill Trewhella (Univ Sydney) 6 members	(2013) <i>Structure</i> 21: 875-881	
Hybrid Methods Task Force	2014	Andrej Sali (UCSF), Torsten Schwede (Univ Basel), Jill Trewhella (Univ Sydney) 27 members	(2015) <i>Structure</i> 23: 1156-1167	



# Presenting Data Quality to Diverse Audiences

- Provide relative and absolute quality metrics in graphical format
- Provide tabulations of key data, refinement statistics, and quality diagnostics
- Assess all macromolecular and ligand structural components
- PDF format reports can be uploaded with manuscript submission to a journal
- Diagnostics also delivered as an XML format data files

## **Overall Quality**



## **Residue Plots**





# **OneDep – The PDB Biocuration Platform**



## **RCSB PDB Data Delivery Pipeline**





# **RCSB PDB Web Portal**

## Launchpad for a wide range of functionalities

Deposit

• Search

Analyze

• Visualize

• Tabulate

## Download

## http://rcsb.org







# **RCSB PDB Mobile App**

- Provides convenient access to PDB data on the go with a minimal feature set
- Provides a browser, simple search, and an interactive 3D viewer
- Supports iPhone, iPad, and Android





http://www.rcsb.org/pdb/static.do?p=mobile/RCSBapp.html



# Web Services (RESTful APIs)

- Programmatic access to data: application-toapplication communication
- Provides external workflows and analysis tools with direct access to a wide range of PDB data and services
- Enables integration of PDB data and services with programs and scripts in a variety of computer languages and computing environments

http://www.rcsb.org/pdb/software/rest.do



# Enabling Data Access Through Integration and Visualization

- Protein Feature View: mapping protein sequence to 3D structure
- Gene View: mapping genome location to 3D structure
- Visualization
  - Browser native visualization tools using an efficient data compression protocol
  - Small molecule electron density and binding interactions
- Data integration resource files provided for download:
  - Correspondences with CCDC ligand structures
  - Sequence cluster data files
  - Phased release of data to support blinded molecular docking tests



## **Protein Sequence Integrated View**





## Genome Sequence Integrated View

Genomic coordinates:	Cytogenetic location: 17q21.31	chr17:41,197,694-41,258,543	Chromosome Location -	C reset view
Genome Assembly GRCh37	<b>.</b> 2			
X Fullscreen				
< 17:41,188,36641,2	e61,407 Q	500kb <b>Q</b>	+ 4	<b>↓ ⊖ ¢; ?</b> >
* Genome 41,190,000	41,200,000 41,210,000	41,220,000 41,230,000	41,240,000	41,250,000 41,2
		, ← 1		≪  ←
× m Genes €	HERE ALGK_A_BRCA1: 4IGK_A_	BRCA1 ×		<- <b>↓</b> <
CBHCAT	Location undefined:41197	710-41222995		
		210 1211 4		
× Repeats				
Conservation			a. Andritt in	



## **Molecular Visualization**

5IRE

The cryo-EM structure of Zika Virus



## http://mmtf.rcsb.org/ https://github.com/arose/ngl

### Large Structure Visualization with NGL/MMTF 08/23

#### Scalable Molecular Visualization

Fast, interactive display of molecular complexes containing millions of atoms is now possible, without plug-ins, on both desktop computers and smartphones. Our new molecular visualization tool, the NGL Viewer, enables easy visualization of even the largest structures for everyone in research and education. In addition, we have developed the new binary compressed Macromolecular Transmission Format (MMTF) to massively reduce network transfer and parsing times. With these advances, we are prepared for upcoming advances in experimental techniques that are already delivering structures of rapidly increasing size and complexity.



Ownload Files -

Display Files

## Visualizing Small Molecule Interactions

Small Molecules Ligands 1 Unique The SMALL ID **MOLECULES** section Chains Name / Formula / InChI Key 2D Diagram & Interactions **3D** Interactions on entry's Structure REA RETINOIC ACID Δ Ligand Explorer Query on REA Summary page C20 H28 O2 SHGAZHPCJJPHSC-Binding Pocket (JSmol) offers tools for ligand Download SDF File ④ YCNIQYBTSA-N visualization Electron Density (JSmol) Download CCD File ④ in 2D and 3D. **POSEVIEW 2D** diagrams BINDING POCKET show a ligand and (JSmol) offers a 3D interacting residues. view of residues contact Black dashed lines with the selected ligand. indicate hydrogen bonds, A cropped ligand van salt bridges, and metal der Waals surface will interactions. Green solid appear that is color line show hydrophobic coded by the proximity interactions and green to the van der Waals dashed lines show surface of the binding  $\pi$ - $\pi$  and  $\pi$ -cation site (red: close contact,

Biotin bound to Streptavidin PDB ID: 1STP



interactions.

Visualization of Biotin bound to Streptavidin PDB ID: 1STP with PoseView (poseview.zbh. uni-hamburg.de)



Electron density mini-map for Retinoic Acid in PDB ID 1CBS at 1sigma contour level. In this example, the ligand fits well into the density.



Maraviroc with its receptor C-C chemokine receptor type 5 PDB ID: 4MBS



LIGAND EXPLORER is a Java Web Start application with options to display hydrogen bonds, hydrophobic contacts, water mediated hydrogen bonds, metal interactions, and surfaces that can be

color-coded using various parameters.



# **Reaching Diverse User Communities**

Who are our users?	What are they using?
<b>Biologists</b> : structural biology, biophysics, biochemistry, genetics, Immunology, pharmacology, cell and molecular biology	RCSB PDB website, deposition tools, data
<b>Other scientists</b> : bioinformatics, software developers,	Web Services, search engines, data
Students & teachers	PDB-101
<b>Media</b> : Writers, textbook authors, patient advocacy groups,	Images, data, information, outreach material, e.g., posters
General public: Curious/interested individuals, artists, sculptors,	Images, Molecule of the Month, information from external media





Online Educational Resources **Resources to help** understand biology at the molecular level

http://pdb101.rcsb.org/

### Animations



## **Paper Models**



Glucose is transported through the blood,

powering cells throughout the body. However, if the level of glucose gets too high, it can cause

the chronic problems associated with diabetes

Atomic structures have revealed how glucose levels are regulated in the body and are providing

To learn more, browse PDB-101 resources

related to diabetes, including features on

new hope for combating the disease.

glucagon, insulin, and RAGE.

PDB-101

Educational portal of

Zika Virus

Zika virus

mellitus.

More

Molecule of the Month

Molecular explorations

through biology and medicine





2016 Video Challenge for High School Students

#### Structural Biology and Diabetes

Create a video that tells the molecular story of health

### Posters





Video Challenge Deadline:

New Insulin and Diabetes

Join RCSB PDB and CCDC

May 29

poster

» 04/26/2016

» 05/10/2016



May 2016

### 🖬 🕊 🖬 🗯 🏟 🗘

Browse resources by category

Health and Disease

## Molecule of the Month







### Quasisymmetry in Icosahedral Viruses

Viruses use quasisymmetry to build large capsids out of many small subunits

In the 1950s, before any virus structures had been determined, Watson and Crick predicted that viruses are symmetrical assemblies of many small subunits. They based this proposal on a puzzling observation: only a small amount of RNA (or DNA) can fit inside a tiny virus, but this RNA needs to encode the sequences for the proteins comprising the virus. This problem is particularly severe for the smallest viruses, such as satellite tobacco necrosis virus (shown here from PDB entry 2buk). Its capsid is just big enough to hold a single strand of RNA that encodes only one protein: the capsid subunit. It's called a "satellite" virus because it can't infect a cell on its own. Instead, it relies on other viruses, which infect the cell at the same time, to provide the necessary viral machinery for its reproduction.

#### Quasisymmetry

Viruses typically need to be larger than these tiny satellite viruses, so that they can encode additional viral proteins. One way they do this is by using *quasisymmetry*. In perfectly icosahedral viruses, the subunits form pentamers that assemble at the vertices of a small icosahedron. In quasisymmetrical viruses, a single type of subunit plays several different roles. They form pentamers that assemble at the vertices, but they also assemble in a slightly different way to form hexamers that fill out the faces of larger icosahedrons. In the early 1960s, Caspar and Klug discovered that there are a few systematic ways to arrange these pentamers and hexamers, based on a *triangulation number* that quantifies the number of slightly different shapes that the subunit needs to adopt in the final capsid.

### **Triangulation Numbers**

Perfectly symmetrical viruses have a triangulation number T=1, because all of the subunits are identical. In satellite tobacco necrosis virus, this is just big enough to hold an RNA strand 1239 nucleotides long. Tomato bushy stunt virus (PDB entry 2tbv) has the next larger triangulation number, T=3, with pentamers on the vertices of the icosahedron and a hexamer centered on each face. It can hold a larger genome with 4776 nucleotides that encodes five proteins, including a viral replicase. The T=4 capsid of Nudaurelia capensis omega virus (PDB entry 1ohf) has hexamers centered on each edge of the icosahedron, and holds two strands of RNA that encode replicase and capsid proteins. Bacteriophage HK97 (PDB entry 1ohg) has T=7 quasisymmetry, and is large enough to enclose a DNA strand with 39,732 nucleotides that encodes 61 proteins.



Icosahedral virus capsids. In each virus, all of the subunits are chemically identical, but they adopt a few different quasisymmetrical shapes, each colored differently here. Pentamers of subunits are colored red, and hexamers of subunits are colored in shades of yellow and orange.

Download high quality TIFF image





PDB members past and present at the PDB40 Anniversary Symposium, 2011

## F Worldwide Protein Data Bank



### pdbj.org



F**unding:** NBDC-JST BMRB bmrb.wisc.edu

Funding: NLM



The Protein Data Bank Archive is managed by:



wwpdb.org

Members



## rcsb.org



Funding: NSF, NIH, DOE BPDBe Protein Data Bank in Europe pdbe.org





Funding: EMBL-EBI, Wellcome Trust, BBSRC, NIGMS, EU



26