# Welcome and Introductions

- RCSB Protein Data Bank Advisory Committee

- Funding Representatives
  - NSF
  - NIH
  - DoE
  - HHMI

- RCSB Protein Data Bank Leadership Team
  - Rutgers, The State University of New Jersey
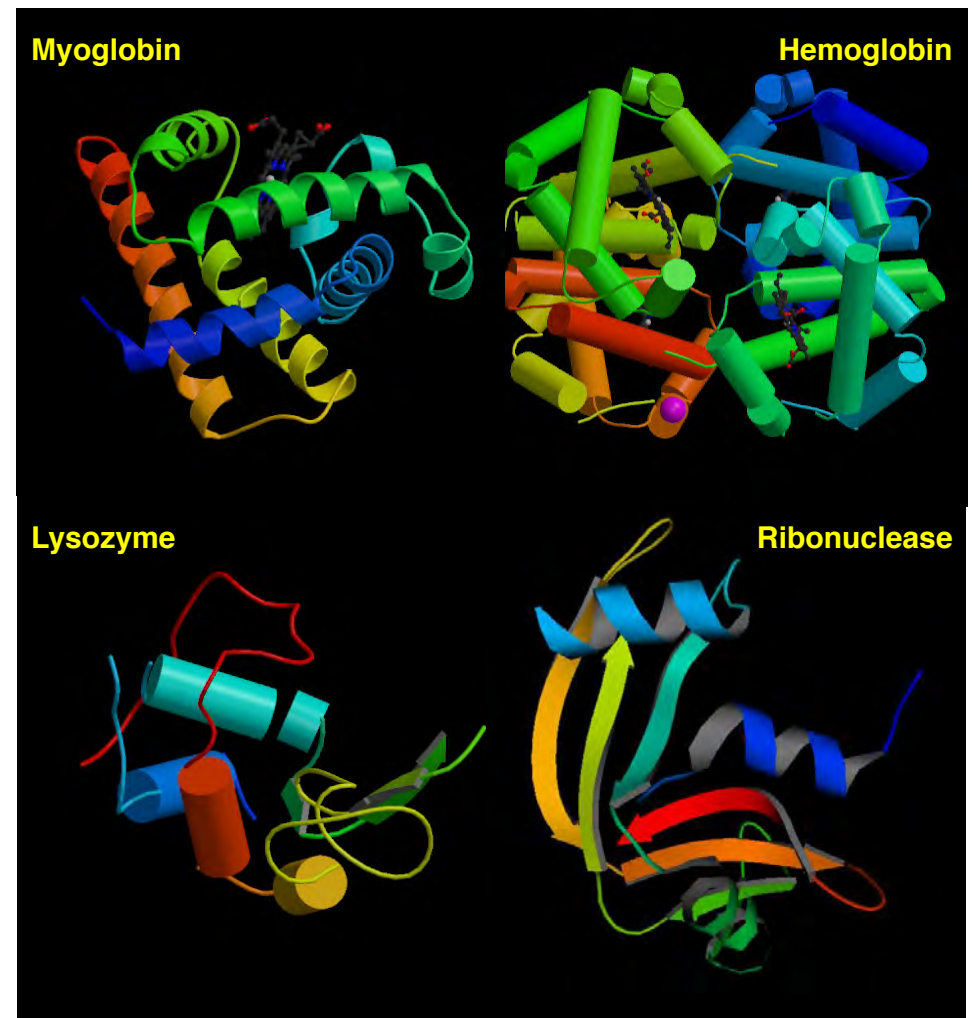  - University of California San Diego

# Protein Data Bank Overview

Stephen K. Burley, M.D., D.Phil.

# Protein Data Bank Archive

- Single primary data archive for 3D structures of proteins, DNA, and RNA

- Established 1971 as 1st Global Open Access digital data resource in biology at Brookhaven (→RCSB PDB 1999)

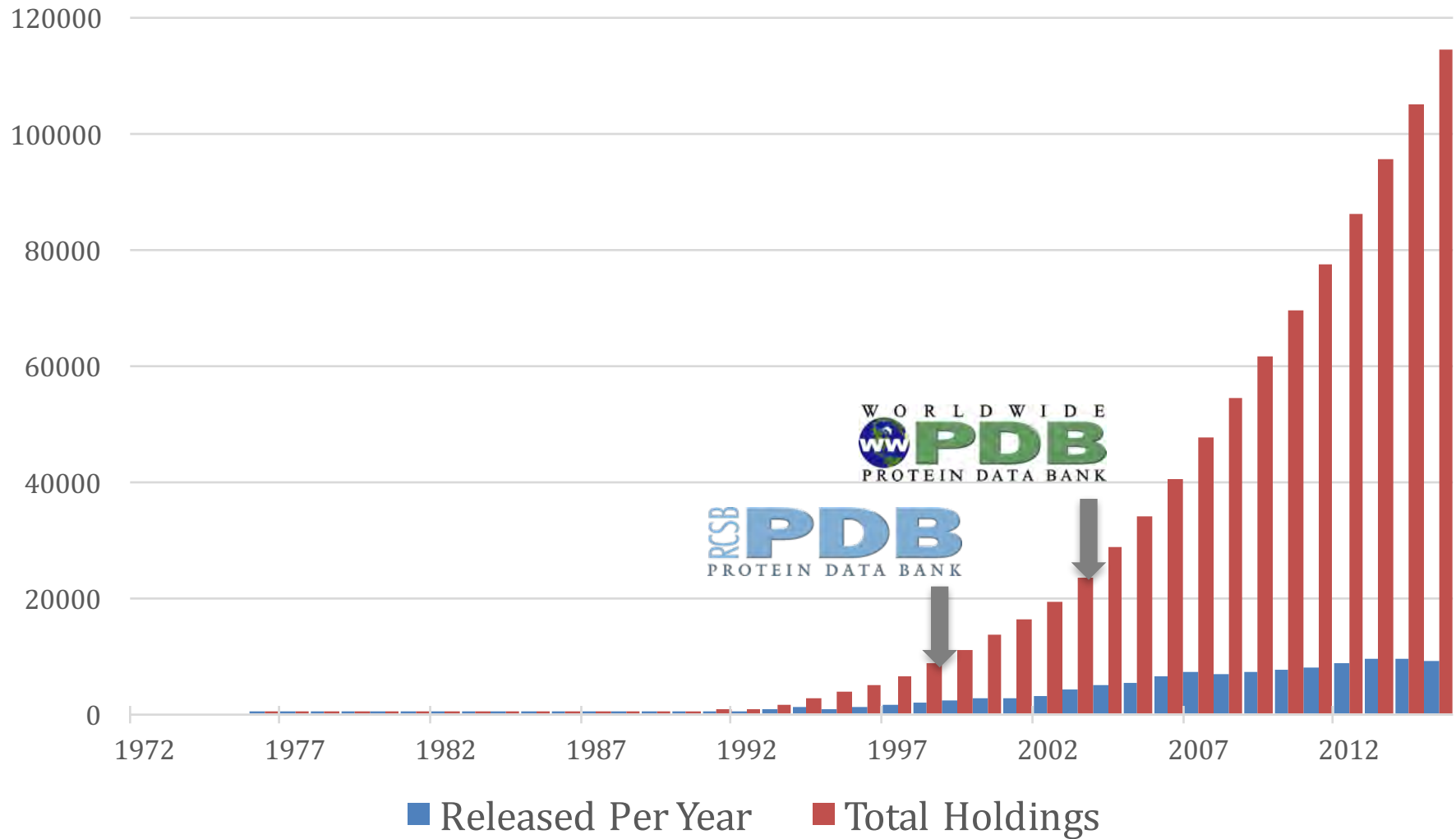- Since 2003, managed by an international partnership (wwPDB)



Myoglobin

Hemoglobin

Lysozyme

Ribonuclease

Some of the very first structures in the PDB
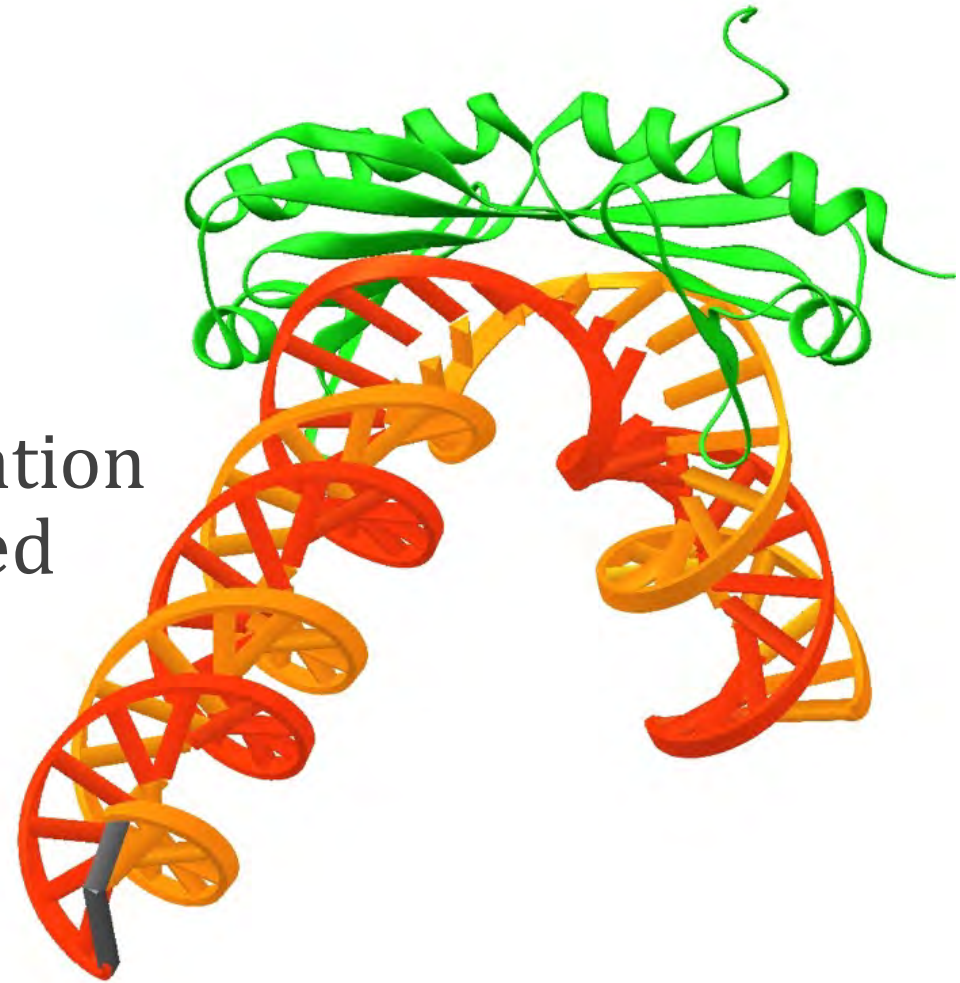
# PDB Commitment to Scientific Data Reuse

- Follow the *FAIR Guiding Principles for scientific data management and stewardship*
  - **F**indability
  - **A**ccessibility
  - **I**nteroperability
  - **R**eusability
- See Wilkinson *et al.* (2016) *Scientific Data* doi: 10.1038/sdata.2016.18

# Released Entries 1971-2015
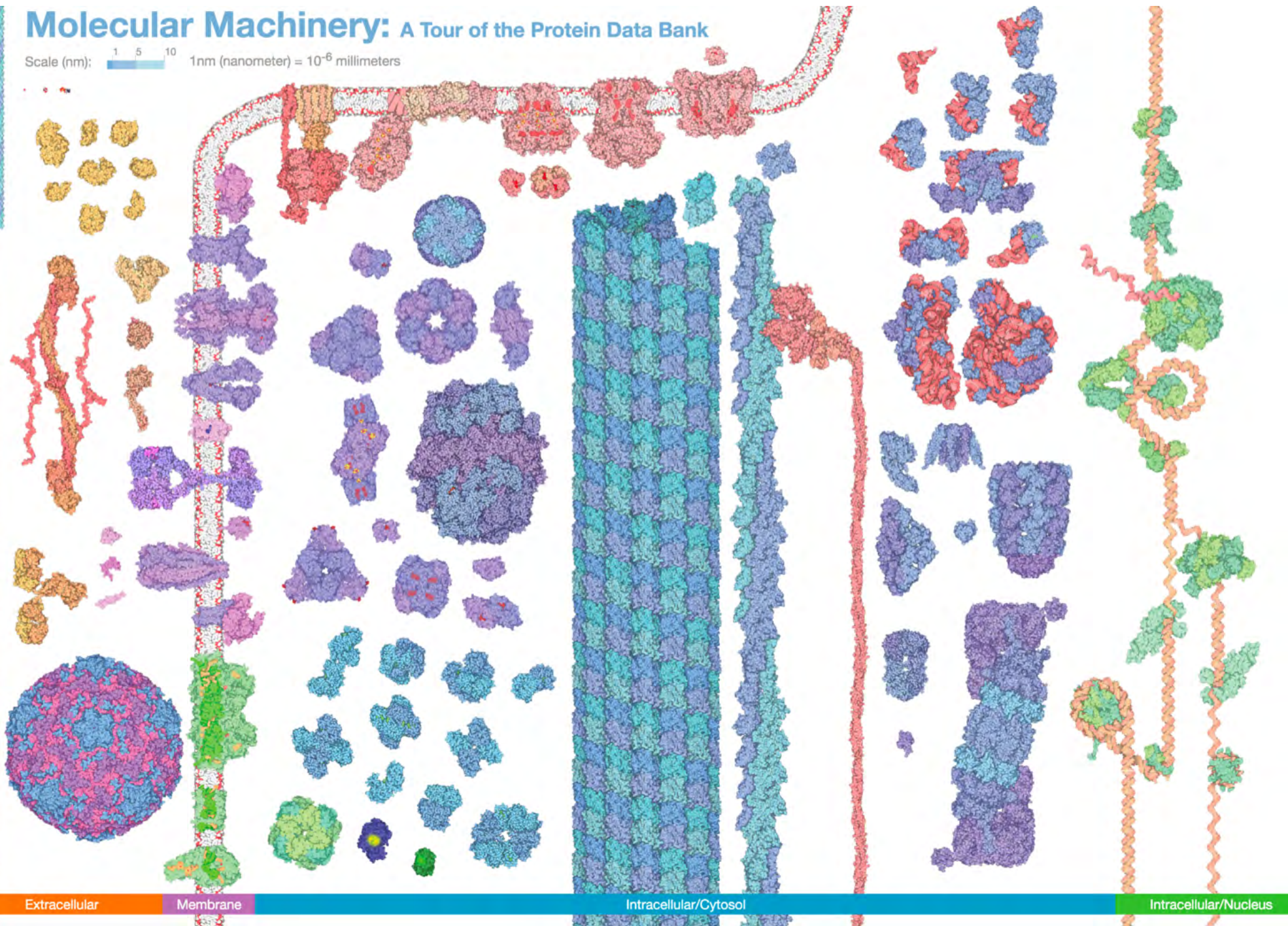
# Function Follows Form in Biology

- 3D structure determines biological/biochemical function

- PDB data inform every area of research and education in biology, basic and applied

- PDB data are used every day to understand health and disease

- PDB data central to drug discovery

*Arabidopsis thaliana*
TATA-box Binding Protein + DNA (PDB 1VTL)

# Molecular Machinery: A Tour of the Protein Data Bank

Scale (nm): 1 5 10    1nm (nanometer) = $10^{-6}$ millimeters
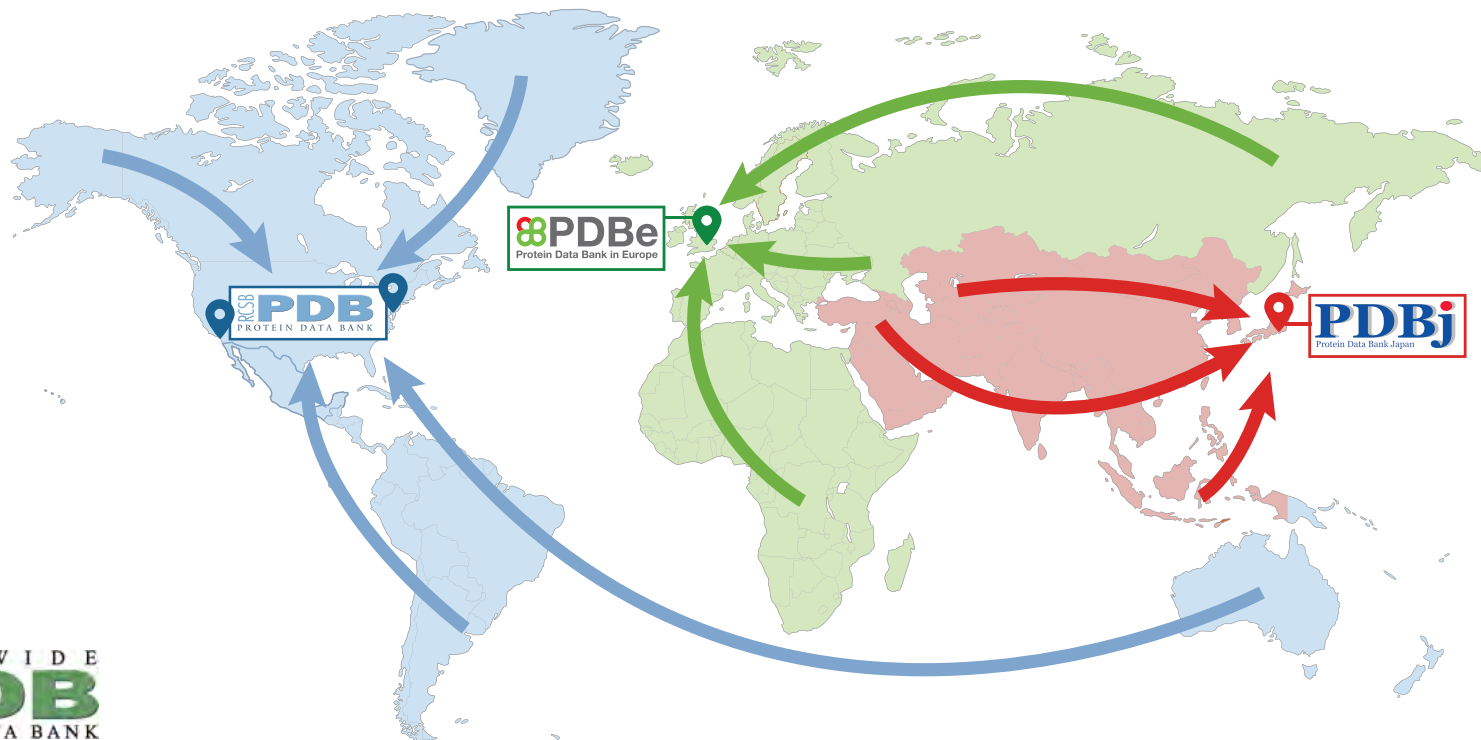
Extracellular    Membrane    Intracellular/Cytosol    Intracellular/Nucleus

# Worldwide Protein Data Bank (wwPDB)

- Structure data are globally produced/consumed

- Regional Data Centers: RCSB PDB (US), PDBj (Asia), PDBe (EU); BioMagResBank (US/Japan)

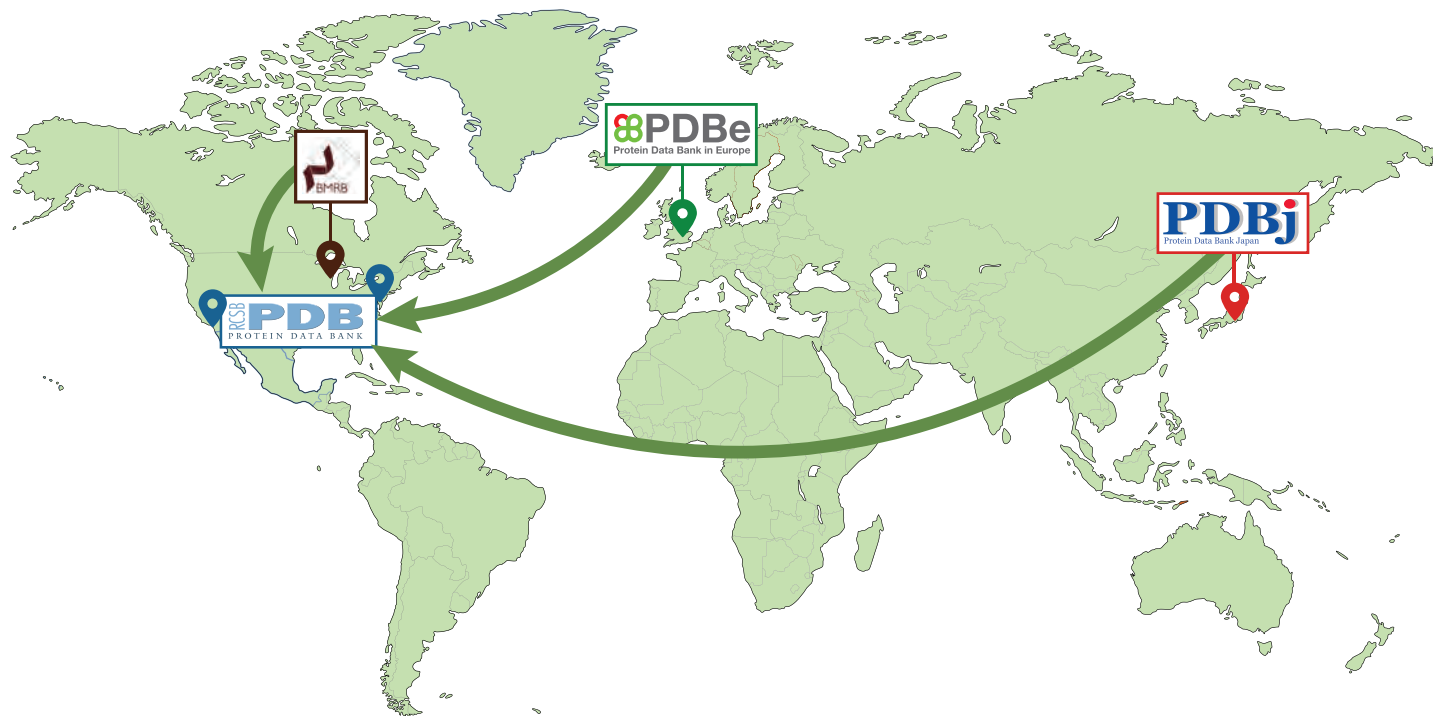- Unfunded; Operated under formal MOU since 2003

# wwPDB Data Centers

- Ensure unrestricted PDB access worldwide

- Work with the scientific community to establish common data standards and best practices

- Collaborate on Global "Data In" Services: Deposition/Biocuration/Validation

- Operate identical FTP data distribution sites

- Develop/provide complementary Global Services for "Data Out"

# RCSB PDB is the PDB Archive Keeper

- Support data security and global disaster recovery

- Ensure data uniformity and consistency

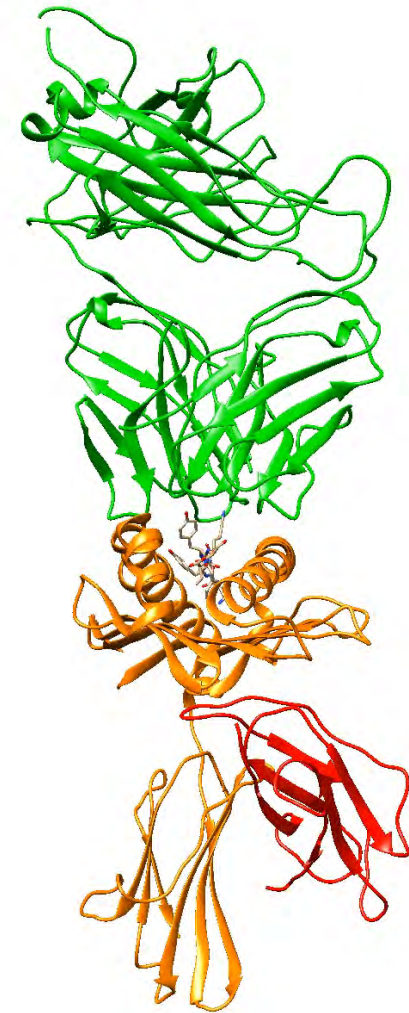- Coordinate weekly updates and FTP distribution

# Cost of Replicating the PDB Archive

- Data integrity and security are of paramount importance to the wwPDB partnership

- Estimated cost of replicating each PDB entry ranges from US$50,000 to > US$250,000

- Cost of replicating the PDB archive **US$12 billion** (assuming <unit cost>=US$100,000)

- Absent PDB data sharing, structural biology would never have reached current heights

# What Has the PDB Archive Enabled?

- Reproducibility and Secure Storage

- Accelerated structure determination technologies

- Understanding evolution in 3D
  - Structure classification and prediction

- Structure-based drug discovery

- Functional understanding of Biology at molecular and atomic levels

Antigen Presenting Cell meets the T-cell
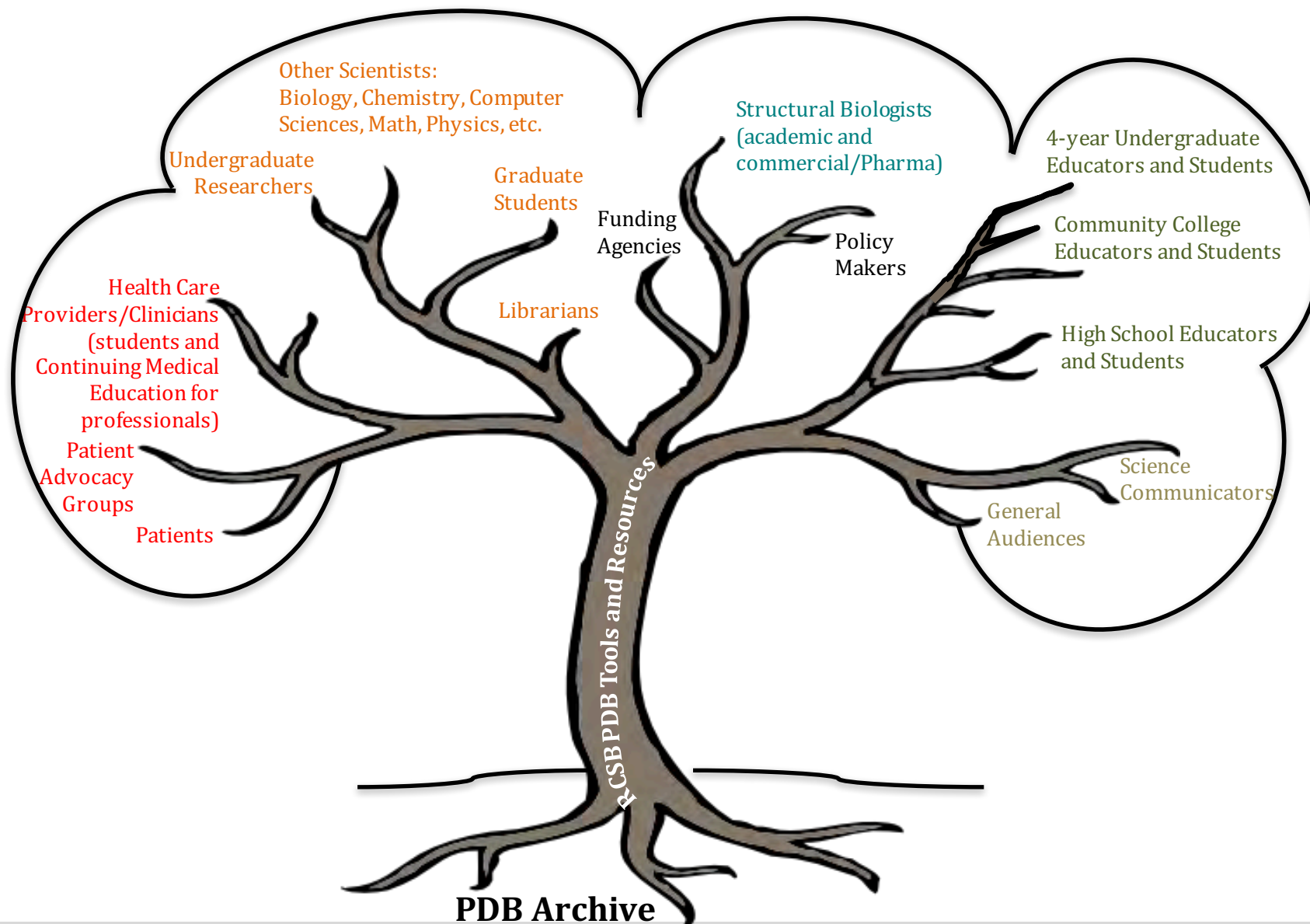PDB 2CKB, Garcia *et al.* (1998)

# PDB Archive Facts and Figures

- Archival Contents
  - ~124,000 Structures Released since 1971
  - ~11,000 New Structures Deposited/Year

- Global User Base
  - ~30,000 Depositors Worldwide
  - >1 Million Unique Visitors/Year
    from 192/195 UN-recognized sovereign nations

- Impacts all of Biology and Medicine
  - >500 Million Data Files Downloaded/Year
  - ~1.5 Million Data Files Downloaded/Day
  - >200 derived data resources repackage PDB data

# RCSB PDB: US wwPDB Data Center

- Established 1999 (RUTGERS | UC San Diego SDSC)

- Founded wwPDB in 2003 to support Data Producers

- Collaborates with international experts and resources to support Data Consumers

- Core activities funded by NSF (DBI-1338415), NIH, DOE

- Competes for additional funding for value-added activities

# RCSB PDB Serves Diverse Stakeholders



Other Scientists: Biology, Chemistry, Computer Sciences, Math, Physics, etc.

Undergraduate Researchers

Graduate Students

Structural Biologists (academic and commercial/Pharma)

4-year Undergraduate Educators and Students

Health Care Providers/Clinicians (students and Continuing Medical Education for professionals)

Funding Agencies

Policy Makers

Community College Educators and Students

Librarians

High School Educators and Students

Patient Advocacy Groups

Science Communicators

Patients

General Audiences

RCSB PDB Tools and Resources

PDB Archive

# RCSB PDB Core Responsibilities

# Archive Keeping

- PDB Archive FTP sites 24/7/366 availability with security and disaster recovery

- Failover and load balancing between Rutgers and UCSD

- Continuous global monitoring



status.rcsb.org

# Data In/wwPDB

- Deployed unified global deposition system (OneDep)

- Validation for the entire PDB archive

- Facile management of very large structures (HIV capsid: 2.4 MM atoms)

### 3J3Q
Atomic-level structure of the entire HIV-1 capsid

Note: Use your mouse to drag, rotate, and zoom in and out of the structure. Click to identify atoms and bonds.

NGL is a WebGL based 3D viewer powered by MMTF.

# FTP Data Distribution/Access



The RCSB PDB Private Cloud

Compute & Storage Nodes

Management Nodes

High Capacity Networks

RCSB PDB distributes data using a Private Cloud→Future Public Cloud

# Data Out

- Optimization of rcsb.org website continues

- Integration of genome sequence →protein sequence →3D structure

- Visualization of pathways, ligands, very large structures, *etc*.



Erlotinib targeting EGFR for Lung Cancer (PDB 4HJO)

# Education

- Developed modular curriculum on diabetes for use in high schools

- Deployed diabetes materials (Molecule of the Month articles, poster) *via* PDB-101 website



Insulin Hexamer/Monomer from Molecule of the Month on *Designer Insulins* (February 2016)

# Global Outreach and Engagement

- Wellcome Image Award successes

- Deployment of Zika virus outreach materials with Purdue structure release

- Publication of Ligand Validation Workshop whitepaper in *Structure*

OVERALL WINNER

Ebola Virus painted for Molecule of the Month by David Goodsell

# RCSB PDB Organization

**Office of Director**

*Stephen K. Burley* (Director, Principal Investigator)*
*Helen M. Berman* (Director Emerita)*
*Christine Zardecki* (Deputy Director)
Luz Fajardo (Administrative Assistant)

## Biocuration

*Jasmine Young**
*Irina Persikova*
*Yuhe Liang*
Luigi Di Costanzo
Sutapa Ghosh,
Brian Hudson, Ezra Peisach
Monica Sekharan
Chenghua Shao, Lihua Tan
Marina Zhuravleva

## Software Development

**East**
*John Westbrook**
*Zukang Feng*
Li Chen, Vladimir Guranović
Rob Lowe, Raul Sala
Wendy Tao, Huanwang Yang

**West**
*Peter Rose*, Andreas Prlić**
Ali Altunkaya, Chunxiao Bi,
Anthony Bradley, Jose Duarte,
Tara Kalro, Jesse Woo

## Systems Administration

**East**
*Harry Namkoong*
Ken Dalenberg

**West**
*Cole Christie*
Chris Randle

## Outreach and Education

*Shuchismita Dutta**
*Christine Zardecki*
David Goodsell
Rachel Green
Maria Voigt

Key:
*Leadership in italics*
* Presenting Today

# RCSB PDB Staff Outreach Commitment

NJ Science Olympiad

San Diego Science and Engineering Festival

ABRCMS 2016

Undergraduate Research,
Summer 2016, Rutgers

High School and Undergraduate Research,
Summer 2016, UCSD

# Supporting Diversity

- Longtime commitment to supporting a diverse, inclusive, and family-friendly workplace

- Mentorship of students under-represented in sciences
  - External *via* Rutgers Office of Diversity and Inclusion Summer Program (RiSE, 3 students in 2016)
  - Rutgers Undergraduates (4 students in 2016)
  - UCSD Outreach Programs

- Joint participation with the Rutgers Center for Graduate Recruitment, Retention, and Diversity at ABRCMS and SACNAS 2016 national meetings

# RCSB PDB Advisory Committee

- Provides independent advice to
  RCSB PDB Director and staff
    - Operates under formal Terms of Reference
    - Triennial rotation schedule (renewable)
    - Cynthia Wolberger agreed to chair through 2019

- Comments, advises, or makes recommendations for
  action on topical issues as they arise over the course of
  the time between meetings, and on any standing
  agenda items
    - Deposition Policies and Annotation Practices
    - Data Distribution, Query Policies, and Practices
    - Education and Outreach

# Agenda

| | |
|---:|:---|
| **Overview** | Stephen K. Burley |
| **Data In** *OneDep, Data Standards, Infrastructure, Plans Forward* | Jasmine Young and John Westbrook |
| **Data Out** *Access and Exploration* | Peter Rose and Andreas Prlić |
| **Outreach** | Helen M. Berman |

**Lunch**

| | |
|---:|:---|
| **Education** | Shuchismita Dutta |
| **Funding and Sustainability Response to 2015 Report** | Stephen K. Burley |

**Matters Arising & General Discussion**

# Data In:
# OneDep, Data Standards, Infrastructure, Plans Forward

Jasmine Young, Ph.D.

John Westbrook, Ph.D.

**rcsb.org**

# Outline

- Team and responsibilities

- Data life cycle

- wwPDB international partnership

- Importance of biocuration

- Engaging scientific communities

- Infrastructure

- Plans forward

# Biocurators and Data In Developers

- 11 scientists,
  3 scientist programmers,
  4 software developers

- 14 Ph.D., 2 M.S., 2 B.S.

- 8 countries,
  3 continents

- Combined length of service 169 years

- Median length of service 9.4 years

# Data In Responsibilities



RCSB PDB curates >5,000 entries/year

Develop and maintain tools for Deposition, Biocuration, and Validation

Assure quality of weekly data release

Make larger archive-wide improvements ("remediation")

Communicate with Depositors and diverse User Communities

Create, maintain and refine data dictionaries

Resource integration across Data In pipeline

Primary Data Curation

Outreach/ Customer Support

wwPDB Deposition, Biocuration (OneDep)

Project Management

Data Standards and New Content

Archive Quality Control

Resource Development

# Data Life Cycle



**wwPDB OneDep**
Unified global deposition, biocuration, and validation system

**3** Deposition
deposit.wwpdb.org

**2** Pre-deposition Validation
validate.wwpdb.org

**4** Biocuration

**Archive Keeper**
RCSB PDB packages and re-distributes data to wwPDB partners

**1** Data Production

**5** Archive Update
ftp://ftp.wwpdb.org

**Structure Biologists**
Generate atomic model and data files

Assemble mandatory data items for deposition

**6** Public Release
http://rcsb.org
ftp.rcsb.org

**Data Users**
Enables research in various fields

Data accessed *via* ftp download or web service

# wwPDB International Collaboration (Data In)

- Developing unified global tools for deposition, validation, and biocuration

- Defining data standards and content: PDBx/mmCIF Dictionary

- Ensuring data uniformity in the PDB archive ("Remediation")

- Maintaining a single global archive

# RCSB PDB is the Archive Keeper

- RCSB PDB is the Archive Keeper for world distribution of PDB data and leads the wwPDB collaboration in developing tools, setting data standards, and performing data remediation

# Unified Global OneDep Tool

- More complete data capture

- File format standardization

- Improved efficiency and consistency
  - Enables workload balancing
  - More automation
  - File replacement pre-submission

- Validation for all methods
  - Standalone Validation Server
  - Web Service API

- Support for larger and more complex structures


2007 Initial Discussions

2010 OneDep Team Meeting



2016 OneDep Summit Meeting

# RCSB PDB Effort on OneDep Project

- Provided technical and managerial leadership

- Responsibilities
  - Backend infrastructure and technology
  - Biocuration pipeline
  - Hosting wwPDB development servers
  - Technical support for partner server installation

- >50% of wwPDB-committed FTEs from RCSB PDB
  - Jasmine Young, Global Project Lead

# Workload Balancing/Depositor Support

deposit.wwpdb.org

Others 19%

USA 81%

45% Americas, Oceania

36% Europe, Africa

19% Asia

# Increasing Size and Complexity

## Number of Large Structures Deposited



Legend:
- Number of large structures (chains > 62 & atoms > 99999)
- Number of structures with MW > 500,000

## Number of Ligands Released





Faustovirus (PDB 5j7v, Klose et al., 2016) is the largest PDB structure





Antibiotic quinupristin/Dalfopristin bound to ribosome (PDB 4u26 Noeske et al., 2014)

# Importance of Biocuration

- Enforces data standardization through policies and common biocuration practices

- Ensures data quality and provides value-added annotation

- Communicates possible errors to Depositors (wwPDB Validation Report)

- Maintains data uniformity and compliance in the PDB archive to enable data search and exploration

- Requires domain expertise
  - Cannot be replaced by purely computational means



2014 wwPDB Biocurator Summit



2015 PDBj Biocuration Training

# Data Quality and Value-Added Annotation

- Consistency checking
  - Polymer sequence and taxonomy
  - Ligand stereochemistry
  - Ligand density fit

- Integration with external data resources

- Overall quantitative and qualitative review of deposited data



REA in PDB 1CBS
(Kleywegt et al., 1994)
RSR=0.10, RSCC=0.96

TMP in PDB 3HW4
(Kaushik et al., 2013)
RSR=0.41, RSCC=0.57

# Improving Data Quality

| Validation Server & API | Deposition | Biocuration | Public Release |
|---|---|---|---|
| Pre-validate data independently before deposition | Mandatory acknowledgement of report produced during deposition | wwPDB-recommended report for journal submission | Report available for all released PDB entries |

Coordinates and data frequently replaced
during Deposition and Biocuration

Validation Report submission during manuscript review process

- Mandatory: *Nature* journals, *Acta D & F, FEBS, JBC, J Immunology, eLIFE,* and *Angew Chem Int Ed Engl*

- Recommended: *Cell, Molecular Cell, Structure*

# Key Features of wwPDB Validation Reports

- Graphical overview of data quality

- Residue plots

- Atomic model quality

- Experimental data quality

Overall Quality



Residue Plots



🟩 OK    🟨 1 issue    🟧 2 issues    🟥 3+ issues

🟦 ill-defined by coordinates

# OneDep Biocuration Processing Time

(a) ~1hr: Simple structures without issues

(b) ~4 hrs: More complex structures without issues

(c) ~15 hrs: Structures with issues, including Depositor response time

Legacy ADIT system: 4-5 days

# OneDep Biocuration Impact

## Top issues frequently raised during Biocuration



Number of Instances (y-axis), Type of Issues (x-axis)

- Chirality Error: 46%
- Polymer Backbone Linkage: 44%
- Atomic Clashes: 22%
- Sequence Discrepancy: 10%
- Unrealistic Occupancy: 10%

New data sets received in response to issues raised during Biocuration

- ~29% of all 2015 entries
- ~25% of all 2016 entries

New tools to promote pre-deposition validation

- Standalone Validation Server (now supports NMR, 3DEM)
- Web Service API

# Improving X-ray Structure Quality

Structure quality improvement since the advent of the wwPDB Validation Report



Rfree    Clash    %Rama    %Rota    %RSRZ

Yellow: Legacy System 2012-2013
Green: OneDep 2014-2015

# Enabling Data Search and Exploration

- Data standardization significantly impacts data query

- RCSB PDB has taken leadership
  - Developing and maintaining data standards
  - Data remediation

- 5 rounds of data remediation for entire PDB archive from 2007 to 2014

- OneDep system also supports data remediation
  - 3DEM (in collaboration with EMDataBank)
  - Carbohydrates

Henrick, *et al.,* 2008, *NAR*; Lawson, *et al.,* 2008, *Acta Cryst. D*; Dutta, *et al.,* 2014, *Biopolymers*

# Depositor/User Feedback

- Daily communication between Biocurators and External Users (Depositors)

- ACA, IUCr meetings: demonstrations, posters and exhibit booth

- Internal Users (Biocurators)

- Continuous testing and improvement

- Weekly cross-site reviews of issues

# Enabling Bulk Depositions from Industry

"Group" Deposition developed to meet community need

- Requirements set by wwPDB OneDep Team
- Support for D3R Blind Challenges
- Depositors: Roche, EMD Serono, University of Marburg, University of Essex
- 364 depositions in single group processed in 5 days

**deposit-group.rcsb.rutgers.edu/groupdeposit/**



RCSB PDB GroupDep System

wwPDB OneDep System

# Engaging Scientific Communities

- Defining data content and quality standards

- Task Forces and Working Groups
  - Validation Task Forces (VTFs)
    - X-ray, NMR, 3DEM
  - Small Angle Scattering
  - Integrative/Hybrid Methods Task Force
  - PDBx/mmCIF Working Groups
  - NEF Working Group

- Ligand Validation Workshop

# Defining Data Content & Quality Standards

| Task Force | Meeting/ Workshop | Chair(s)/Membership | Outcome |
|---|---|---|---|
| X-ray Validation | 2008 2015 | Randy Read (Univ of Cambridge) 17 members | (2011) *Structure* 19: 1395-1412 |
| NMR Validation | 2009, 2011, 2013 (x2), 2015 2016 | Gaetano Montelione (Rutgers) Michael Nilges (Institut Pasteur) 10 members | (2013) *Structure*, 21: 1563-1570 |
| 3DEM Validation | 2010 | Richard Henderson (MRC-LMB) Andrej Sali (UCSF) 21 members | (2012) *Structure* 20: 205-214 |
| Small-Angle Scattering | 2012, 2014 | Jill Trewhella (Univ Sydney) 6 members | (2013) *Structure* 21: 875-881 |
| Hybrid Methods | 2014 | Andrej Sali (UCSF), Torsten Schwede (Univ Basel), Jill Trewhella (Univ Sydney) 27 members | (2015) *Structure* 23: 1156-1167 |



W O R L D W I D E
**PDB**
P R O T E I N   D A T A   B A N K

# Community Data Standards

- Data managed using PDBx/mmCIF
  - Extends earlier IUCr data standard
  - PDBx/mmCIF dictionary has >4400 data terms



mmcif.wwpdb.org

- Extensions now coordinated with wwPDB PDBx/mmCIF Working Group
  - Supports broader needs of both contributors and users of the archive

- Host community workshops

- mmCIF.wwpdb.org provides data dictionaries, schema, software tools

# Data Standards Working Groups

- ## PDBx/mmCIF Working Group
  - Experts and methods developers
  - Ensures good support in key community software tools

- ## NMR Task Force Working Group
  - NMR Exchange Format

- ## SASCIF
  - PDBx-compatible extension dictionary supporting data exchange with SASBDB

PDBx Workshop, October 2014

NMR Workshop, August 2016

nature structural & molecular biology
Home | Current issue | Comment | Research | Archive ▾
home ▸ current issue ▸ correspondence ▸ full text

NATURE STRUCTURAL & MOLECULAR BIOLOGY | CORRESPONDENCE    OPEN

NMR Exchange Format: a unified and open standard for representation of NMR restraint data

Appl Cryst JAC  Journal of Applied Crystallography    IUCr IT WDC

search IUCr Journals    GO

home    archive    editors    for authors    for readers    submit    subscribe    open access

JAC CIF APPLICATIONS

*J. Appl. Cryst.* (2016). **49**, 302-310
doi:10.1107/S1600576715024942

OPEN ACCESS

Extension of the sasCIF format and its applications for data processing and deposition

M. Kachala, J. Westbrook and D. Svergun

# Focus on Ligand Quality

Current summary validation statistics may not identify poor electron density fit

NADP in PDB 1ZK4:
2|Fo|-|Fc| map at 1σ
Schlieben *et al.,* 2005

NADP in PDB 2FZD:
2|Fo|-|Fc| map at 1σ
Steuber *et al.,* 2006



| Metric | Percentile Ranks | Value |
|---|---|---|
| Rfree | | 0.175 |
| Clashscore | | 7 |
| Ramachandran outliers | | 0 |
| Sidechain outliers | | 1.0% |
| RSRZ outliers | | 2.4% |

Worse — Better
■ Percentile relative to all X-ray structures
▯ Percentile relative to X-ray structures of similar resolution

| Metric | Percentile Ranks | Value |
|---|---|---|
| Rfree | | 0.169 |
| Clashscore | | 5 |
| Ramachandran outliers | | 0 |
| Sidechain outliers | | 3.2% |
| RSRZ outliers | | 5.1% |

Worse — Better
■ Percentile relative to all X-ray structures
▯ Percentile relative to X-ray structures of similar resolution

**Overall Ligand Model Quality**: Legacy 2012-2013 (yellow) *vs*. OneDep 2014-2015 (green)



RSR   RSCC   Bond RMS   Angle RMS   OWAB   LLDF   RSZD+   RSZD-

# Ligand Validation Workshop

- Co-crystal structure determination experts (Academe and Industry) and Software Developers (X-ray Crystallography and Computational Chemistry) discussed, developed, and recommended:
  - Best practices for PDB archive deposition/validation of co-crystal structures
  - Editorial/Refereeing/Publication standards for co-crystal structures



Adams *et al.* (2016) *Structure* 24: 502-508

Rutgers July 30-31, 2015

# High Availability OneDep

- Redirection between global partners sites in the of event loss-of-service

- RCSB PDB hosts bi-coastal OneDep services
  - Warm failover:  Independent East and West Coast OneDep installations
  - Active failover:  In-progress East Coast deposition sessions mirrored to West Coast



RCSB PDB West          RCSB PDB East

# Continuous Monitoring

- PDB archive FTP

- wwPDB OneDep systems

- wwPDB validation servers

- wwPDB website

- RCSB PDB website and related services



status.rcsb.org

# Data In Plans Forward



Improve automation in Biocuration and data remediation

Continue to engage Depositors and diverse User communities

Harden OneDep infrastructure to support new wwPDB members (*e.g.,* China and India)

Primary Data Curation

Outreach/ Customer Support

wwPDB Deposition, Biocuration (OneDep)

Project Management

Data Standards and New Content

Archive Quality Control

Resource Development

Continue data standardization efforts to support data remediation

Streamline weekly OneDep update operations

File versioning to improve archive management

Resource integration across Data In pipeline

# Data Out:
# Access, Exploration, and Metrics

Peter Rose, Ph.D.

Andreas Prlić, Ph.D.

# Outline

- Team and responsibilities

- Who uses RCSB PDB?

- What type of research do we enable?

- How broad is our impact?

- Plans forward

# Bi-coastal Developer Team

- 3 scientists-software dev., 6 software developers, 2 systems & infrastructure

- 4 Ph.D., 5 M.S., 2 B.S.

- 9 countries/3 continents

- 48 combined years of service

- 3 median years of service

# Responsibilities



Website and resources utilized by >350,000 unique visitors/month >1 million unique visitors/year

Communicate with diverse User communities

Outreach/ Customer Support

Web Resources

Data Integration and Visualization

>40 External Resources

Project Management

Ongoing usage assessment to inform development

Analytics

Hardware

System Design, Procurement, Maintenance

Archive Keeping

24/7/366 Support, Weekly Updates, Annual Snapshots

# Who Are Our Users And How Do We Know?



rcsb.org



pdb101.rcsb.org

Google Analytics - Server Logs - External Metrics

# RCSB PDB Traffic Patterns/Usage

This is what a busy day looks like.
Arrows indicate highly trafficked user patterns.
Number of page views per day shown.



**RCSB.org Home page** 20,000

**Search** 40,000

**Structure Summary Pages** 50,000

**Data Download**

**PDB-101** 1,400

**Molecule of the Month** 6,500

**3D Visualization** 7,100

**Other**

**Landing pages**

**Programmatic Access**

**FTP/rsync** 1,000,000 downloads/day

**REST API** 120,500 hits/day

5

# Searching and Search Results

- Text search with autosuggestion

- Advanced Search

- New search features
  - Synonyms
  - Sequence clusters
  - PDB-101 content

- New Search Results page
  - Improved usability
  - Responsive layout
  - Fast page loads

# Structure Summary Page

- Entry at a glance

- Detailed data organized in tabs

- New features
  - Reorganized content
  - Responsive layout
  - Integrated data views
    - Validation
    - Web-friendly NGL 3D viewer
    - Protein Feature View
    - Gene View
    - Pathway View

# Impact of Design Changes

Growth in number of page views after October 2015 redesign

# Downloading and Reporting

- Tools used manually and programmatically

- Download structures, ligands, sequences, experimental data

- Report creation for search results

- New Features
  - User Interface overhaul
  - Sequence cluster information
  - Programmatic access improvements (RESTful web services)

# Data Integration Enabling Research I

## Protein Feature View

- Sequence level view of protein features from UniProt, Pfam, PDB, Protein Model Portal, ...

- New tracks
  - SNPs
  - Mutations in PDB
  - Protein modifications
  - Exon mapping
  - Link sequence position to 3D visualization



EGFR mutation
PDB 5HG8

Integrating Genomic Information with Protein Sequence and 3D Atomic Level Structure at the RCSB Protein Data Bank
*Bioinformatics* 2016 doi:10.1093/bioinformatics/btw547

10

# Data Integration Enabling Research II

**Pathway View**

- Maps structures and metabolites to pathways

- Preview on Structure Summary page

- Interactive browsing

- Pathway name searching



Blue: structures in PDB
Yellow: homology models

# Data Integration Enabling Research III

## Gene View

- Mapping structural coverage onto human genome

## Map Genomic Position

- To chromosome, protein position, and 3D structure

# Data Integration Enabling Research IV

## Electron Density

- Mini-maps to assess Ligand quality



Good fit: REA in PDB 1CBS



Bad fit: PDB 3HW4

## NGL Validation 3D View

- Protein structure *versus* Electron density (Real Space R-factor Z score)



Good fit: PDB 3WY6



Bad fit: PDB 3WYZ

13

# Enabling Drug Discovery in 3D



**Ligand Summary**

**Structure Summary**



Gefitinib binding site view, PDB 4WKQ



Gefitinib electron density map

# Visualizing Large Structures in 3D

- 68 of the 100 largest PDB structures were deposited in the past three years

- Challenge: Data transmission and parsing
  - Developed MacroMolecular Transmission Format (MMTF)
  - Rapid adoption by community
    - Jmol, 3Dmol.js, iCn3D(NCBI), PyMol
    - BioJava, Biopython

- Challenge: Web visualization
  - NGL Viewer efficiently renders large complexes using MMTF on any device



**PDB archive parsing time comparison**

MMTF Time, min: **2**

mmCIF     MMTF



NGL is a WebGL based 3D viewer powered by MMTF.

HIV Capsid PDB 3J3Q, ~2.4MM atoms

MMTF
mmtf.rcsb.org

15

# Factors Influencing Decision Making



| 2013-2018 NSF Proposal | User Feedback | Advisory Committees |

**RCSB PDB Data Access and Exploration**

| Data Archive Content | Usage Analysis | Expert Staff Knowledge |

# Website Usage

## Users

- >350,000 monthly
  >1 million annually
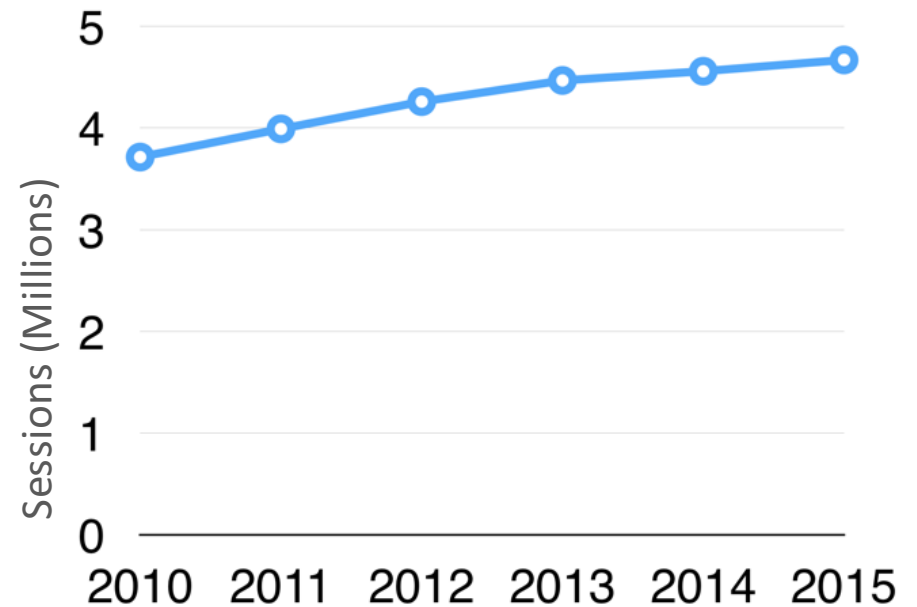
- 2% annual growth in users
  (non-bounce sessions)

## Sessions

- ~26% growth since 2010

- In 2015, ~1 million more
  sessions than in 2010
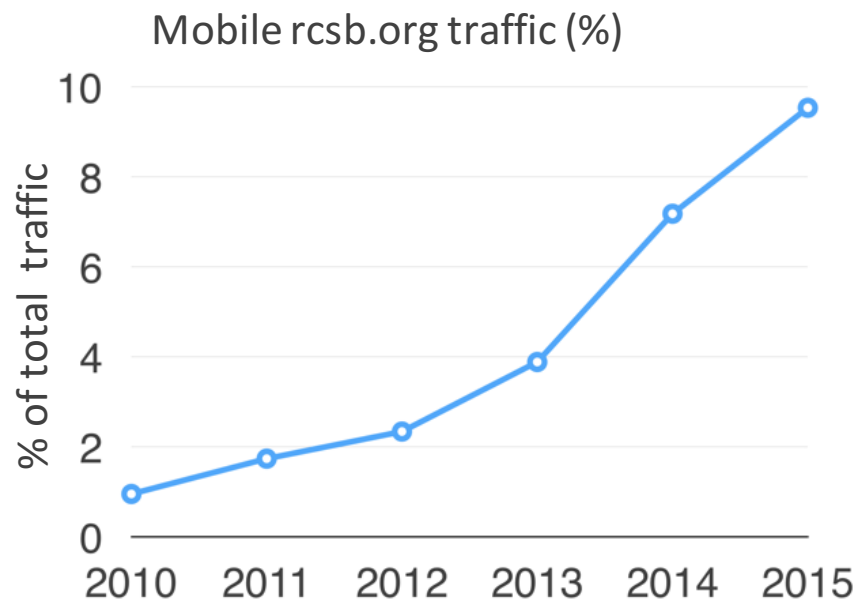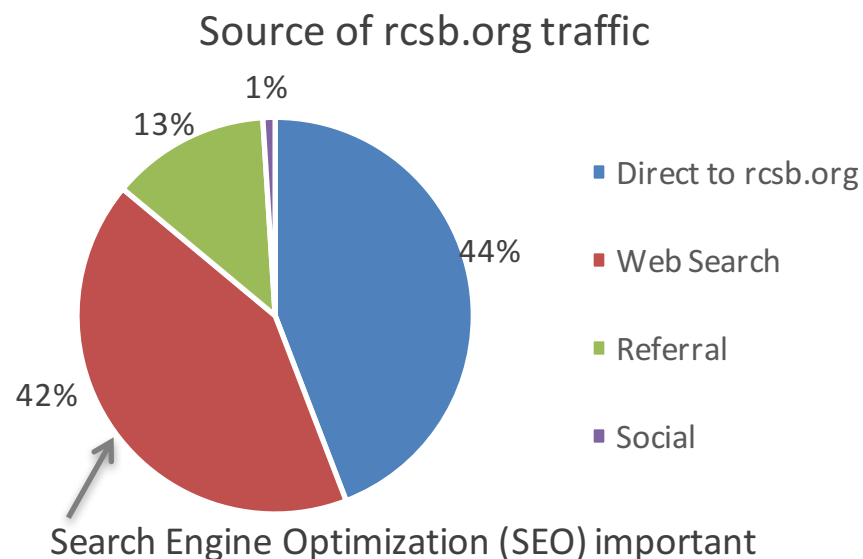
- High average session
  duration (~8 minutes)

Global Traffic

**29%**
US traffic

Non-bounce Sessions

# Website Access

- Most traffic: direct access and web searches (*e.g.*, Google)

- Mobile usage growing rapidly
  - 10% to rcsb.org
  - 20% to PDB-101 (educational site)

- Supported by responsive layout and WebGL-based 3D visualization

Source of rcsb.org traffic



- Direct to rcsb.org
- Web Search
- Referral
- Social

44%

42%

13%

1%

Search Engine Optimization (SEO) important

Mobile rcsb.org traffic (%)



% of total traffic

2010 2011 2012 2013 2014 2015

18

# Data Downloading/Programmatic Access

**Access *via* Website and FTP**

**Programmatic Access (API)**

Entry Downloads and Views 2015

RESTful Web Services:
fastest growing service

# Impact: Primary RCSB PDB Publication

Cited by 21459

Cited ~1500 times/year

Google Scholar

## The Protein Data Bank

Helen M. Berman[1,2,*], John Westbrook[1,2], Zukang Feng[1,2], Gary Gilliland[1,3], T. N. Bhat[1,3], Helge Weissig[1,4], Ilya N. Shindyalov[4] and Philip E. Bourne[1,4,5,6]

[1]Research Collaboratory for Structural Bioinformatics (RCSB), [2]Department of Chemistry, Rutgers University, 610 Taylor Road, Piscataway, NJ 08854-8087, USA, [3]National Institute of Standards and Technology, Route 270, Quince Orchard Road, Gaithersburg, MD 20899, USA, [4]San Diego Supercomputer Center, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0505, USA, [5]Department of Pharmacology, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0500, USA and [6]The Burnham Institute, 10901 North Torrey Pines Road, La Jolla, CA 92037, USA

NEWS FEATURE

## THE TOP 100 PAPERS

Nature explores the most-cited research of all time.

BY RICHARD VAN NOORDEN, BRENDAN MAHER AND REGINA NUZZO

nature MOST CITED PAPERS

PDB #92

| Field: Research Areas | Record Count | % of 15137 | Bar Chart |
|---|---|---|---|
| BIOCHEMISTRY MOLECULAR BIOLOGY | 7907 | 52.236 % | |
| CHEMISTRY | 3075 | 20.314 % | |
| BIOPHYSICS | 2823 | 18.650 % | |
| COMPUTER SCIENCE | 2310 | 15.261 % | |
| PHARMACOLOGY PHARMACY | 1962 | 12.962 % | |
| MATHEMATICAL COMPUTATIONAL BIOLOGY | 1596 | 10.544 % | |
| BIOTECHNOLOGY APPLIED MICROBIOLOGY | 1258 | 8.311 % | |
| SCIENCE TECHNOLOGY OTHER TOPICS | 810 | 5.351 % | |
| CRYSTALLOGRAPHY | 762 | 5.034 % | |
| PHYSICS | 695 | 4.591 % | |
| MATHEMATICS | 648 | 4.281 % | |
| CELL BIOLOGY | 609 | 4.023 % | |
| GENETICS HEREDITY | 390 | 2.576 % | |
| ENGINEERING | 371 | 2.451 % | |
| LIFE SCIENCES BIOMEDICINE OTHER TOPICS | 240 | 1.586 % | |
| MICROBIOLOGY | 177 | 1.169 % | |
| IMMUNOLOGY | 166 | 1.097 % | |
| MATERIALS SCIENCE | 159 | 1.050 % | |
| RESEARCH EXPERIMENTAL MEDICINE | 120 | 0.793 % | |
| PLANT SCIENCES | 118 | 0.780 % | |
| SPECTROSCOPY | 112 | 0.740 % | |
| POLYMER SCIENCE | 96 | 0.634 % | |

WEB OF SCIENCE

2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016

20

# Impact: All RCSB PDB Publications

- 116 publications 2000→

- h-index: ~38

- i10-index: ~62

- Aggregate citations: ~2400/year



| Citation indices | All | Since 2011 |
|---|---|---|
| Citations | 30614 | 13969 |
| h-index | 38 | 33 |
| i10-index | 62 | 53 |



2008 2009 2010 2011 2012 2013 2014 2015 2016

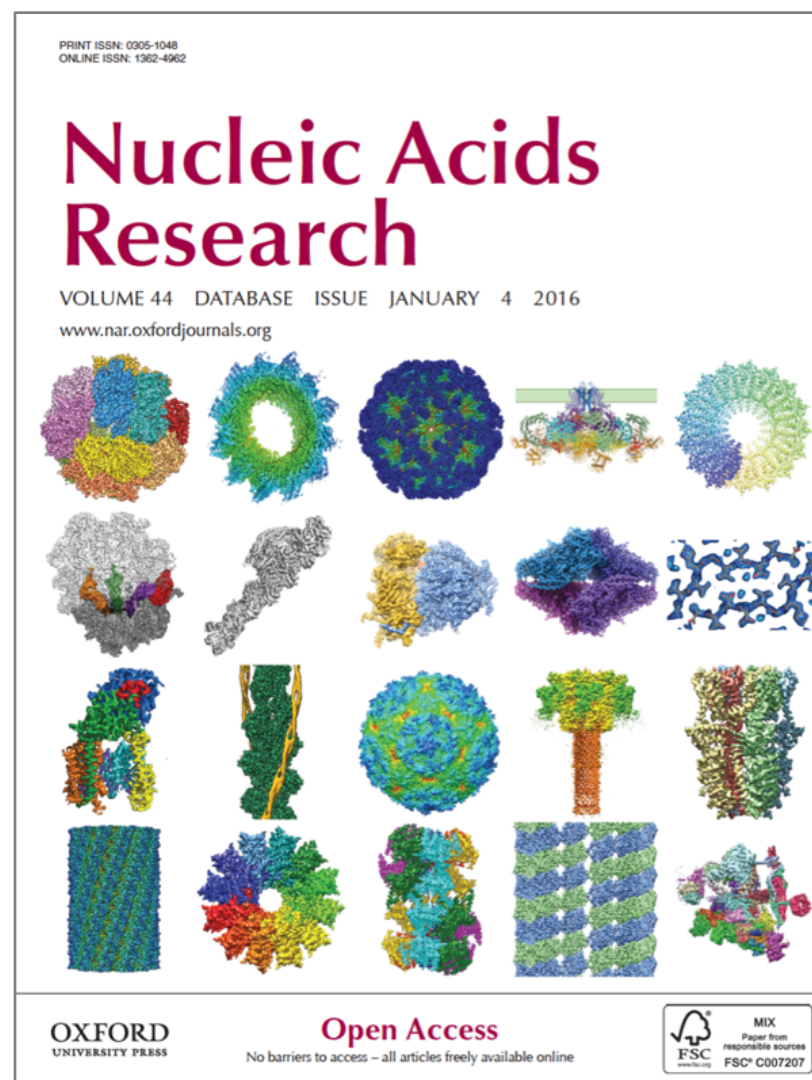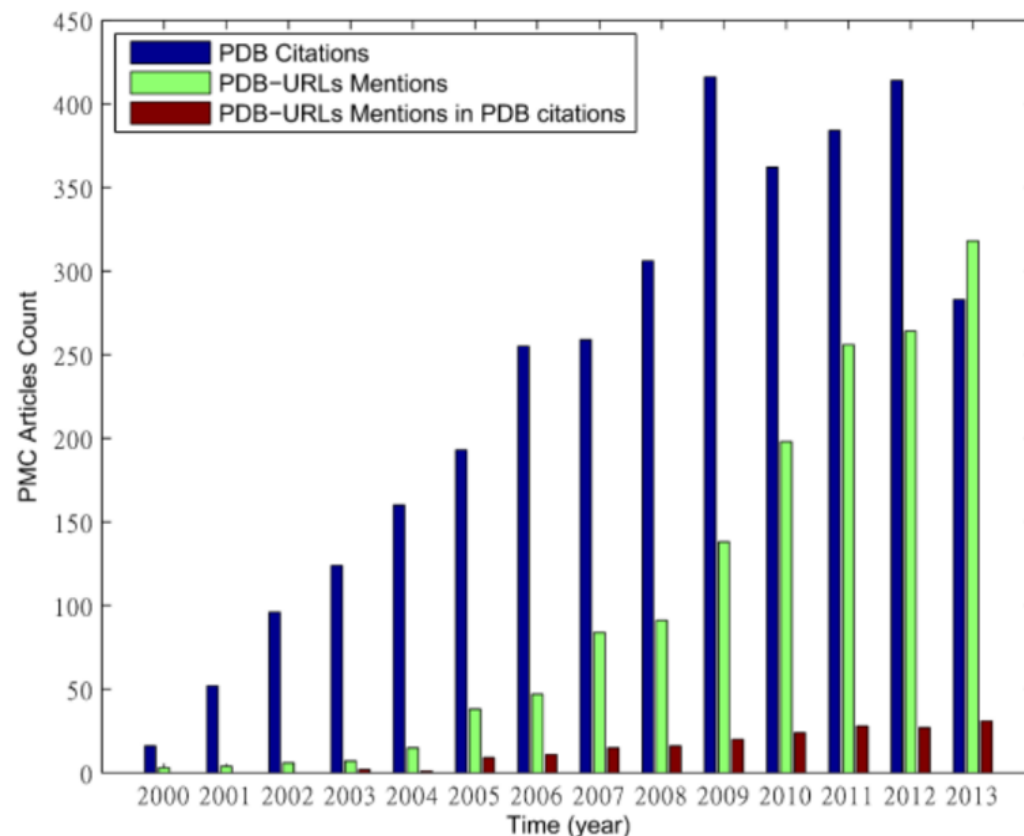| The RCSB Protein Data Bank: redesigned web site and web services PW Rose, B Beran, C Bi, WF Bluhm, D Dimitropoulos, DS Goodsell, ... Nucleic acids research 39 (suppl 1), D392-D401 | 384 | 2011 |
|---|---|---|
| The RCSB Protein Data Bank: new resources for research and education PW Rose, C Bi, WF Bluhm, CH Christie, D Dimitropoulos, S Dutta, ... Nucleic Acids Research 41 (D1), D475-D482 | 274 | 2013 |
| The RCSB Protein Data Bank: views of structural biology for basic and applied research and education PW Rose, A Prlić, C Bi, WF Bluhm, CH Christie, S Dutta, RK Green, ... Nucleic acids research 43 (D1), D345-D356 | 127 | 2015 |

# Impact: PDB Data Reuse

- PDB data used by >200 biological databases
  - Based on databases publishing in *NAR* 2011-2016
  - 11 Categories: Structure, Protein Sequence, Nucleotide Sequence, RNA Sequence, Genomics, Metabolic and Signaling, Human Genes and Diseases, Immunology, Proteomics, Plant, Other

- Since 2011, >25% of new databases utilize PDB data (119 out of 452 new databases)

# Citations in PMC Open Access Articles

- Articles either cite the original PDB publication (Berman NAR 2000) or mention URL rcsb.org
  - Rarely are both referenced

- URL mentions are rising rapidly as data source references

- Citation statistics significantly underestimate the impact of the PDB data resource



Citing a Data Repository: A Case Study of the Protein Data Bank (2015) *PLoS ONE* 10(8): e0136631 doi:10.1371/journal.pone.0136631

# 3166 Patents Mention "protein data bank"

1. 9,476,035 T Recombinant polymerases with increased phototolerance
2. 9,475,886 T Recombinant antibody composition
3. 9,475,881 T Antibody variants with enhanced complement activity
4. 9,475,862 T Neutralizing GP41 antibodies and their use
5. 9,475,851 T High MAST2-affinity polypeptides and uses thereof
6. 9,475,847 T Insecticidal proteins and methods for their use
7. 9,474,759 T Broad-spectrum antivirals against 3C or 3C-like proteases of picornavirus-like supercluster: picornaviruses, caliciviruses and coronaviruses
8. 9,469,684 T Therapeutic and diagnostic cloned MHC-unrestricted receptor specific for the MUC1 tumor associated antigen
9. 9,468,660 T Antinematodal methods and compositions
10. 9,464,311 T Method for identifying modulators of ubiquitin ligases
11. 9,464,280 T Beta-lactamases with improved properties for therapy
12. 9,458,470 T Recombinant influenza virus-like particles (VLPs) produced in transgenic plants expressing hemagglutinin
13. 9,458,434 T Mutant enzyme and application thereof
14. 9,458,229 T Immunogenic proteins and compositions
15. 9,453,236 T Polynucleotides and polypeptides involved in post-transcriptional gene silencing
16. 9,453,224 T MiRNA modulators of thermogenesis
17. 9,453,019 T Linked purine pterin HPPK inhibitors useful as antibacterial agents
18. 9,452,222 T Nucleic acids encoding modified relaxin polypeptides
19. 9,452,210 T Influenza virus-like particles (VLPS) comprising hemagglutinin produced within a plant
20. 9,451,783 T Phytase variants
21. 9,447,157 T Nitration shielding peptides and methods of use thereof
22. 9,447,156 T Methods and compositions for inhibiting neddylation of proteins
23. 9,447,127 T Synthetic lung surfactant and use thereof
24. 9,446,121 T Cloning of honey bee allergen
25. 9,446,116 T Peptide sequences and compositions
26. 9,443,017 T System and method for displaying search results

USPTO PATENT FULL-TEXT AND IMAGE DATABASE

[Home] [Quick] [Advanced] [Pat Num] [Help]

[Next List] [Bottom] [View Cart]

Searching US Patent Collection...

**Results of Search in US Patent Collection db for:**
"protein data bank": 3166 patents.
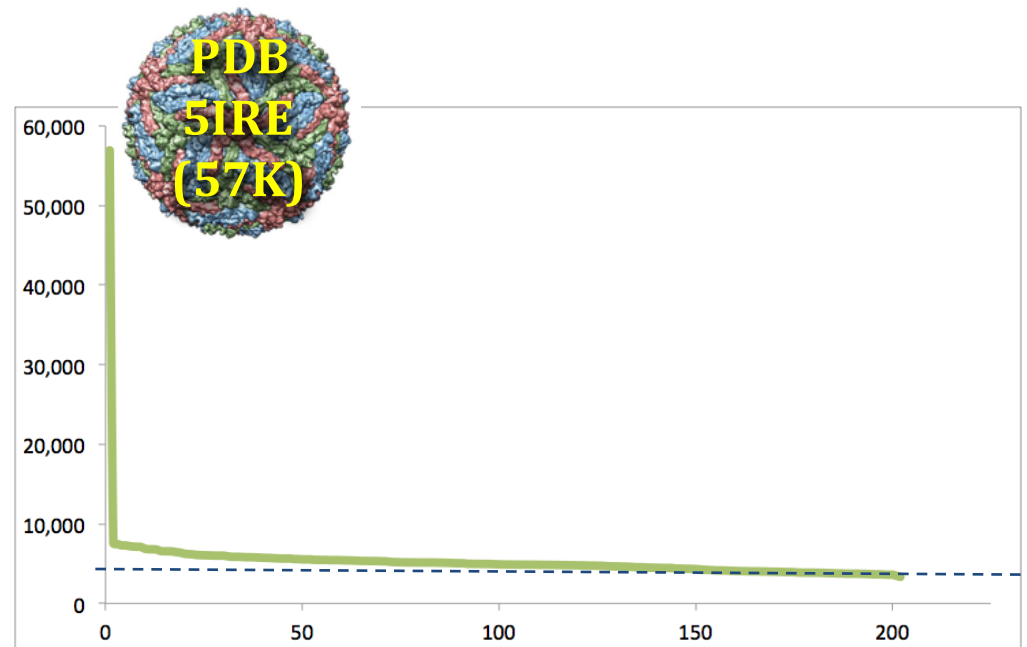Hits 1 through 50 out of 3166

http://patft.uspto.gov/
Accessed October 26, 2016

# Case Study: Zika Virus Data Release

- Zika virus structure PDB 5IRE released March 30, 2016 (Sirohi et al., 2016)

- Downloaded (8K) and viewed on website (49K) times

- >10x usage *versus* 201 other entries released same week
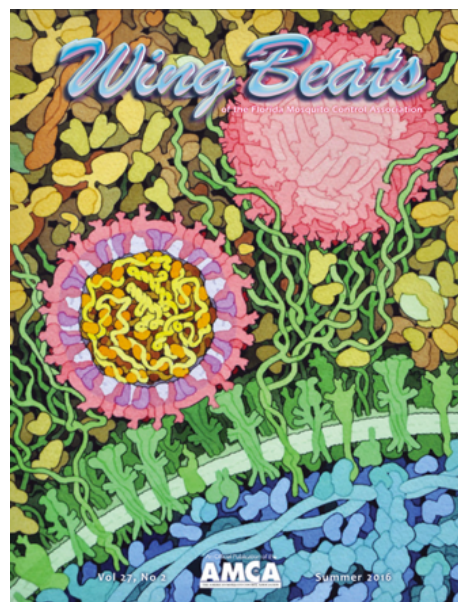  - Average ~5K

Data Downloads & Web Views
March 30-October 4, 2016



202 structures released March 30, 2016, sorted by usage frequency

# Related Outreach

- Molecule of the Month May 2016
  - Zika watercolor highlighted on Cover of *Cell Host & Microbe* and many blogs (NPR, NIH Director, Smithsonian, *Anthropology News*, …)
  - ~12,000 page views

- Molecular Origami PDF paper models

- RCSB PDB Coloring Book

# Infrastructure

- Hosted at SDSC/UCSD and Rutgers

- Disaster preparedness

- Geographic load balancing

- Private cloud
  - Expanded capacity
  - More flexibility → on demand resourcing
  - Better analytics

- High accessibility

# Future Needs

**Continue to transform PDB Data → Knowledge for growing and diverse User base:**

- Data Consumers → expand breadth
  - Easy access and understanding of data
  - Customized views
- Power Users → expand depth
  - 3D queries and mining of PDB
  - Web services (API)

# Data Out Plans Forward

Protein family views, 3D structural queries, APIs

Communicate with diverse User communities, develop reusable open source components

Map genomic (variation) data to structure, link to homology model resources

Web Resources

Outreach/ Customer Support

Data Integration and Visualization

Project Management

Analytics

Hardware

Archive Keeping

Assess usage to inform development (e.g., A/B Testing)

Expand cloud infrastructure to accommodate growing user needs and new features

Support 24/7/366 Operations, Weekly Updates, Annual Snapshots

# Education

## Shuchismita Dutta, Ph.D.

**rcsb.org**

# RCSB PDB User Communities

# Theme-Based Education Strategy

## RCSB PDB Educational Design Process

**In collaboration with**

| RCSB PDB Educational Design Process | In collaboration with |
|---|---|
| Develop/Teach Undergraduate Honors Course | Subject matter-experts |
| Develop online RCSB PDB curriculum | Subject matter-experts<br>High school teachers |
| Test/Refine curriculum | High school teachers<br>Educational consultants |
| Promote curricular modules | High school teachers<br>Educational coordinators<br>Related societies |
| Repackage materials for other audiences | Healthcare professionals<br>Patient advocates |

# Offering Courses and Developing Curricula

# Syllabus of UG Honors Course on Diabetes

## Syllabus (Spring 2015-7)

- Introduction to Insulin and Diabetes

- Understanding the subject matter in 3D

- Clinical aspects of Diabetes and its treatments – Expert lectures

- Approaches to managing Diabetes
  - Non-pharmacological
  - Pharmacological

## Student Assessment Projects

- Molecules involved in glucose homeostasis and causes of Diabetes

- Current pharmacological approaches for treatment of Type 2 Diabetes

Dr. L. Amorosa
RWJMS, Endocrinology

Cynthia Seidman
RD, CDN, Formerly at Rockefeller University

Dr. T. Schneider
RWJMS, Endocrinology

# PDB-101 Resources/Activities



Designer Insulins

MotM, Spring 2016

| Insulin | Feb. 2001 |
|---|---|
| Insulin Receptor | Feb. 2015 |
| Glucagon | Apr. 2015 |
| Receptor for Advanced Glycation End Products | Jun. 2015 |
| Designer Insulins | Feb. 2016 |
| Dipeptidyl Peptidase 4 | Oct. 2016 |



Molecular Origami: Insulin Paper Model and Visualization Activity Fall 2015

# Diabetes Poster



Developed in Spring 2016

# Video Challenge for HS Students

- **2016**
  - Topic: Structural Biology & Diabetes
  - Participation: 82 entries (up from 38 in 2015)
  - Judges:
    - Endocrinology Chief
    - Science Animator
    - Drug Discovery Scientist
    - Scientist/Educator

- **2017**
  - Topic: Treating Diabetes



2016 Judge's Awards First Place

# Offering Courses and Developing Curricula

# Develop/Test Diabetes HS Curriculum

- **Development workshop**
  (July 2016)
  - Draft Modular Curriculum
  - Meet Next Generation Science Standards (NGSS)

- **Recruitment workshop**
  (September 2016)
  - 45 NJ HS Teachers and Science Supervisors attended
  - 28 teachers from 17 schools committed to Pilot Testing

- **NJ Science Convention**
  (October 2016)
  - More recruitment of Pilot Testers

- **Pilot Testing**
  - 2016-17 Academic year



Dr. L. Amorosa
RWJMS, Endocrinology

Dr. A. Ohri
RWJMS, Endocrinology

Dr. M. Kamienski
Rutgers School of Nursing



K. Shah

Standing (L to R): Mr. R. Tempsick, Mr. B. Buck, Ms. S. Coletta, Mrs. A. Sanelli, Ms. J. Jiang, Mrs. H. Sharif, Mrs. S. Eswaran
Sitting (L to R): Dr. B. Ameer, Ms. M. Dominguez, Dr. S. Dutta, Dr. M. Battacharya

# Pilot Testing Diabetes HS Curriculum

## Curriculum at a Glance

- Modular Components
  - Introduction to Proteins
  - Learning to use RCSB PDB data, tools, and resources
  - Enzymes
  - Protein Synthesis
  - Endocrine System
  - Cell Signaling
  - Genetics
  - Evolution
  - Managing Diabetes

- Module content
  - Pre- and Post-Tests
  - Learning Materials with Notes
  - Activities with Teacher Notes

## Where are our Testers?



**28 teachers**
**17 Schools**

# Selected HS Activities



PDB-101 — Molecular explorations through biology and medicine

Educational portal of PDB

Non-public version of the PDB-101 Web site for Diabetes Curriculum Pilot Testing

**Curriculum Modules** — Overview · Discussion Forum · Contact Us · Teacher Log In

Diabetes at a Molecular Level

- Getting Started
- Overall Learning Objectives
- Learning Materials
- Hands-On Activities
- Monitoring Student Learning

Suggestions for Teachers · Skills

## Getting Started

If you are a teacher and would like to access the accompanying teaching notes, click here or use the 'Teacher Log In' link on the Curriculum Modules menu bar.

Before | After

Tree

| A0A0B5AC95 | A0A0B5AC95_CONGE |
| O73727 | INS_DANRE |
| P01308 | INS_HUMAN |
| P01317 | INS_BOVIN |

Weaponized
Insulin

PDB 5JYQ
PDB 1TRZ

Engineered
Insulin
Molecules

K29 Acylated
via spacer

Degludec (Acylated Insulin),
PDB 4AJX

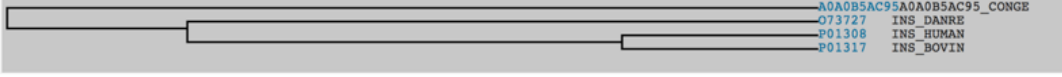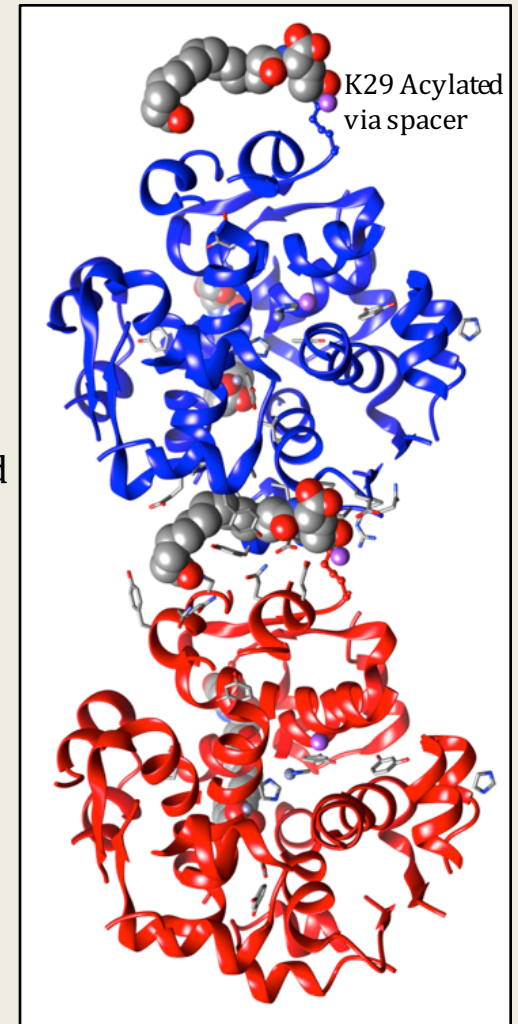# Evaluation and Expansion

## Expert Advisors/Plans

- Karen Collias
  - Founder
- Jennifer Childress
  - Director Instructional Support
- Sue Coletta
  - WSSP/Colleague
- Planned evaluation
  - EQuIP applied to NGSS
    - Alignment to the NGSS
    - Instructional Supports
    - Monitoring Student Progress

## Evaluation Process

- Gather Data: Pilot Testing of Diabetes curriculum
  - Pre- and Post-tests
    - Student Progress
  - Teacher Surveys
    - Alignment to NGSS
    - Instructional Support
  - Student Artifacts
    - Student Progress
- Contract for Professional Evaluation of Curriculum (using above data)

# Summary: Education Efforts

# Funding and Sustainability

# Response to 2015 Advisory Committee Report

Stephen K. Burley, M.D., D.Phil.

# Current Funding

- Core mission support DBI-1338415 for 2014-2018
  (NSF, NIH, DOE; competing renewal likely in 2018)

- Non-core activity support
  - NIGMS Drug Design Data Resource (Amaro/Burley, UCSD)
  - NCI BD2K-Structural Biology Data Compression (Rose, UCSD)
  - NIH BD2K-BioCaddie (Rose, UCSD)
  - NLM BD2K Data Science Course (Lawson, Rutgers)
  - NSF-Integrative/Hybrid Methods EAGER (Berman, Rutgers)
  - NSF-Data Management EAGER (Berman, Rutgers)
  - NSF Big Data Spoke Planning Grant (Prlić, UCSD)
  - NSF REU Minority Summer Students (Burley, Rutgers)
  - NIDA Science Olympiad (Herman/Dutta, Milwaukee School of Engineering)

- RCSB collaborative projects:
  EMDataBank, BioSync, NDB, SBKB

- Private support for Outreach projects
  - HIV film (Viiv, IBM, Rutgers, *et al.*)
  - Symposium on Aesthetics and the Life Sciences
    (Wellcome Trust, Princeton, Rutgers)

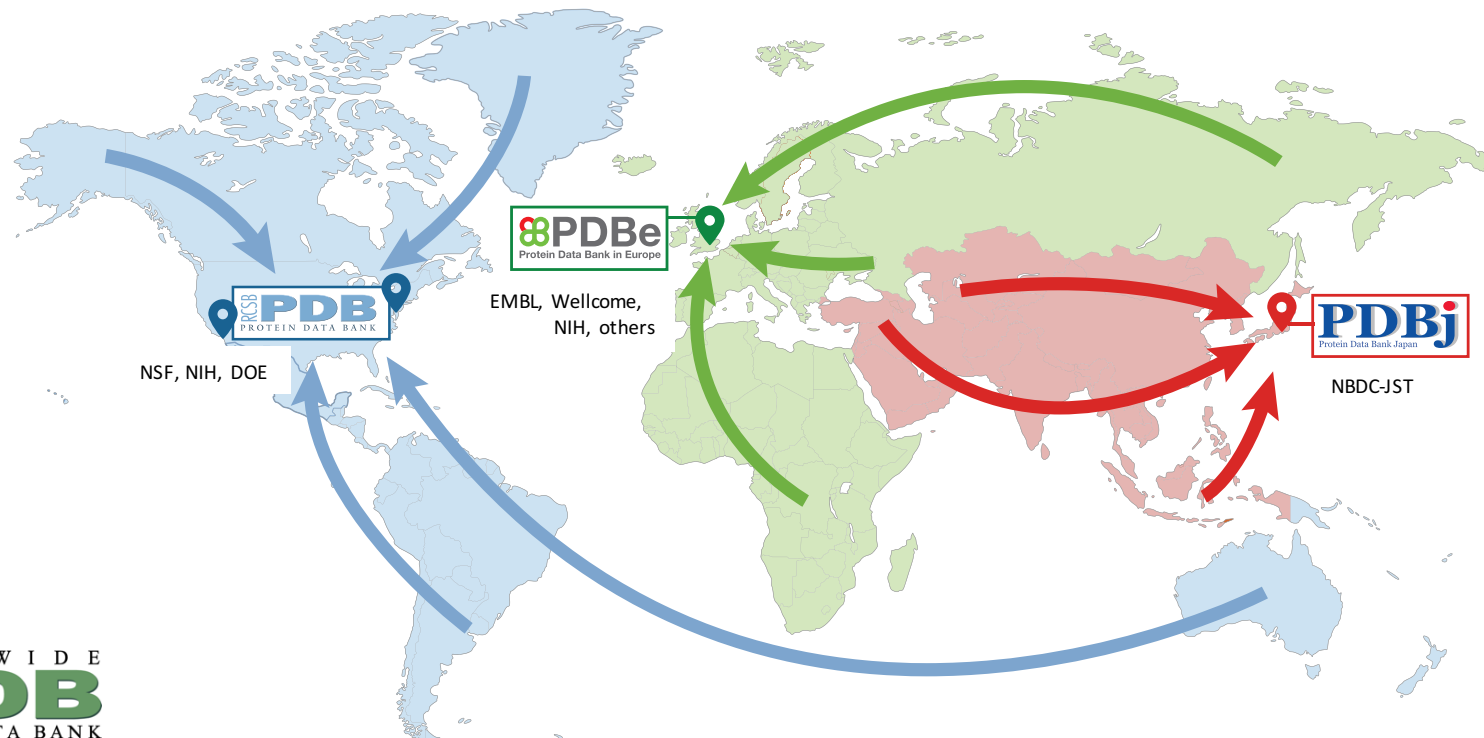# RCSB PDB Funding

- >10 years of nearly flat funding has resulted in a substantial decline in our purchasing power
  - 2004 funding was $5,926,617
    - Equivalent to ~$7,574,649 in 2016 (inflation)
  - 2016 funding is $6,455,369
    - Purchasing power down by ~$1,119K (⬇~14.8%)

- To add "Insult to Injury"
  - 2013 funding was $6,688,486
  - 2016 funding is $6,455,369 (⬇3.5%)
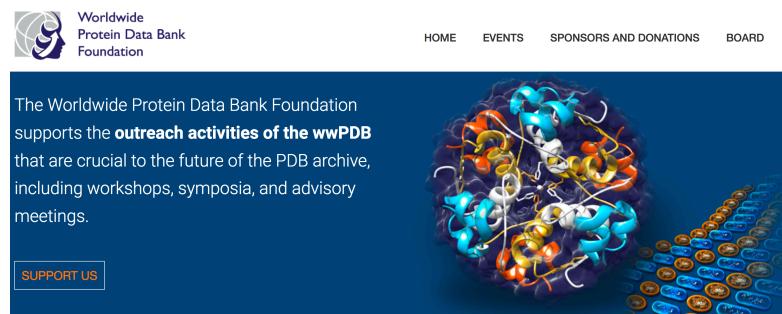
# Data In Sustainability

Worldwide Protein Data Bank (wwPDB)

- Data In shared among Regional Data Centers: RCSB PDB (US), PDBj (Asia), PDBe (EU); BioMagResBank (US/Japan)

# wwPDB Foundation: Outreach Fundraising

- **Established to support specific wwPDB activities**
  - Outreach and education activities, including seminars and workshops
  - Partial support for Advisory Committee meetings
  - Inaugural event: PDB 40th anniversary in 2011
  - Next major milestone: PDB50

- **501(c)3 organization**
  - American, tax-exempt association dedicated to scientific, literary, charitable, and educational purposes

- **Fundraising on-going**



Members of the PDB, past and present, in attendance at PDB40

# Data Out Sustainability

- RCSB.org one of most heavily-used primary biology data resources worldwide
  - Users: >1 million unique visitors/year
  - Global Reach: ~30% US, ~70% non-US
  - Most Data Consumers are not Data Producers

- Core activities enhanced *via* peer-reviewed grant applications for discrete technology development

- Joint wwPDB proposals planned for developing features common to both Data In and Data Out

# Contributions to Sustainability Dialogue

- Sustaining Domain Repositories for Digital Data Working Group (Helen M. Berman)

- Sustaining Biological Infrastructure Advisory Board (Helen M. Berman)

- CODATA/SciDataCon (R. Andrew Byrd, Economics and Impact of the Protein Data Bank (PDB) Archive)

- International Human Frontier Science Program Organization (HFSPO) Life sciences data resources and the future

- NSF Advisory Committee for Cyberinfrastructure (Helen M. Berman)

- Gateways 2016: 11th Gateway Computing Environments Conference



*Longstanding participation in formal and informal sustainability discussions*

# Sustaining Domain Repositories for Digital Data Working Group *Principles*

1.  Research data are a Public Good

2.  Science requires a durable and permanent record

3.  Repositories provide essential domain expertise

4.  Data should be prepared for curation prior to publication (not after the fact or never!)

5.  Sufficient and long-term financial support is critical

6.  Global partnerships, both public and private, should be encouraged

7.  Fiscal transparency is essential

# Sustaining Domain Repositories for Digital Data Working Group *Funding Requirements*

- Economic Stability/Long-term Sustainability

- Global Open Access

- Equity for Data Depositors

- Equity for Research/Teaching Institutions

*8 different funding models examined*
*Only one meets all requirements*

# The Infrastructure Funding Model

- Funding agencies commit to direct payment of the costs of archiving experimental data/metadata generated with the research support they provide

- Data Resource funding comes in the form of strategic, long-term infrastructure investments (divorced from typical 3-5 year grant cycles)

- Ensures Economic Stability/Sustainability for an Open Access Data Resource Ecosystem with Equity for Data Depositors and Consumers

- PDB replacement cost: $US 12 billion

- Estimated archiving cost/year: ~2% of structural biology funding

# Questions for the Committee

- Sustainability
  - Are there other funding sources we should be exploring?
  - Are we making convincing arguments?
  - Are there other data you want to see?
  - Are we reaching the right people?