# Some Economic Considerations for Managing a Centralized Archive of Raw Diffraction Data
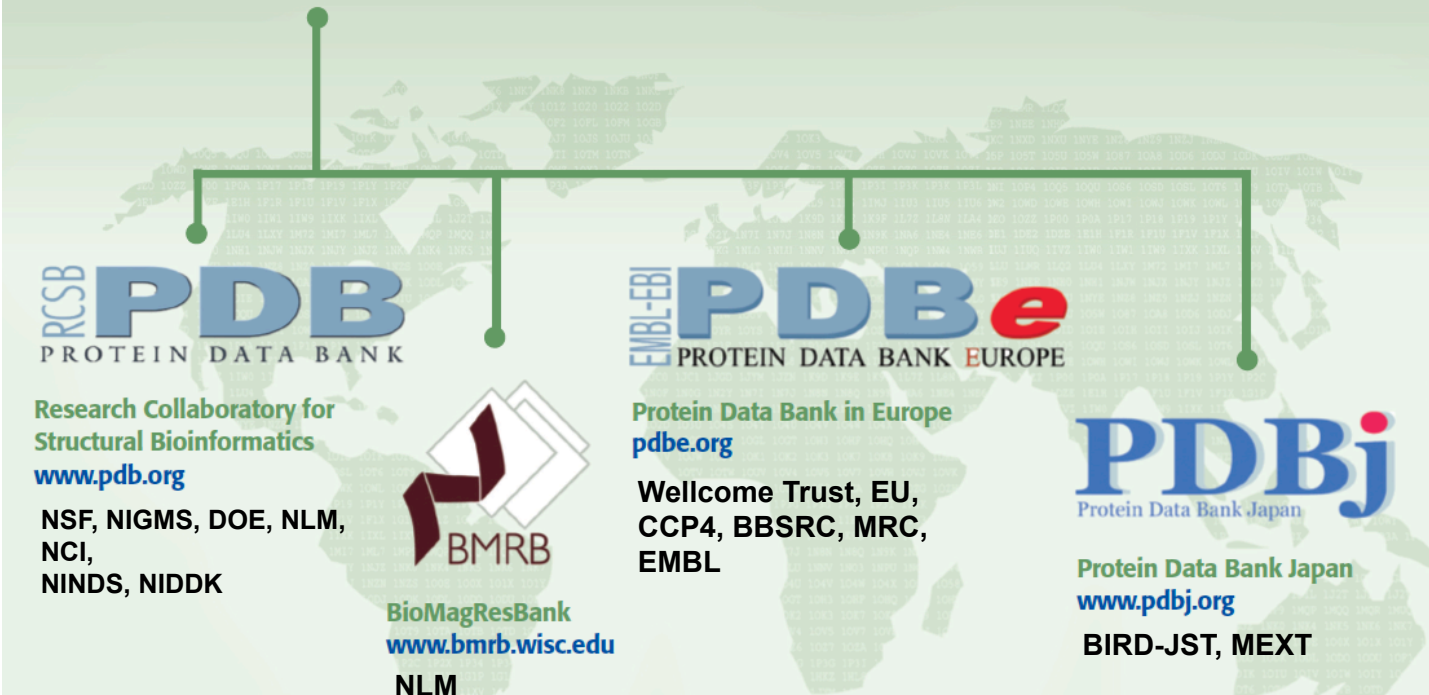
John Westbrook

WORLDWIDE
wwPDB
PROTEIN DATA BANK

www.wwpdb.org

# **Overview**

- PDB as a community partner
- Challenges and scope of archiving primary data
- Some technical and cost alternatives
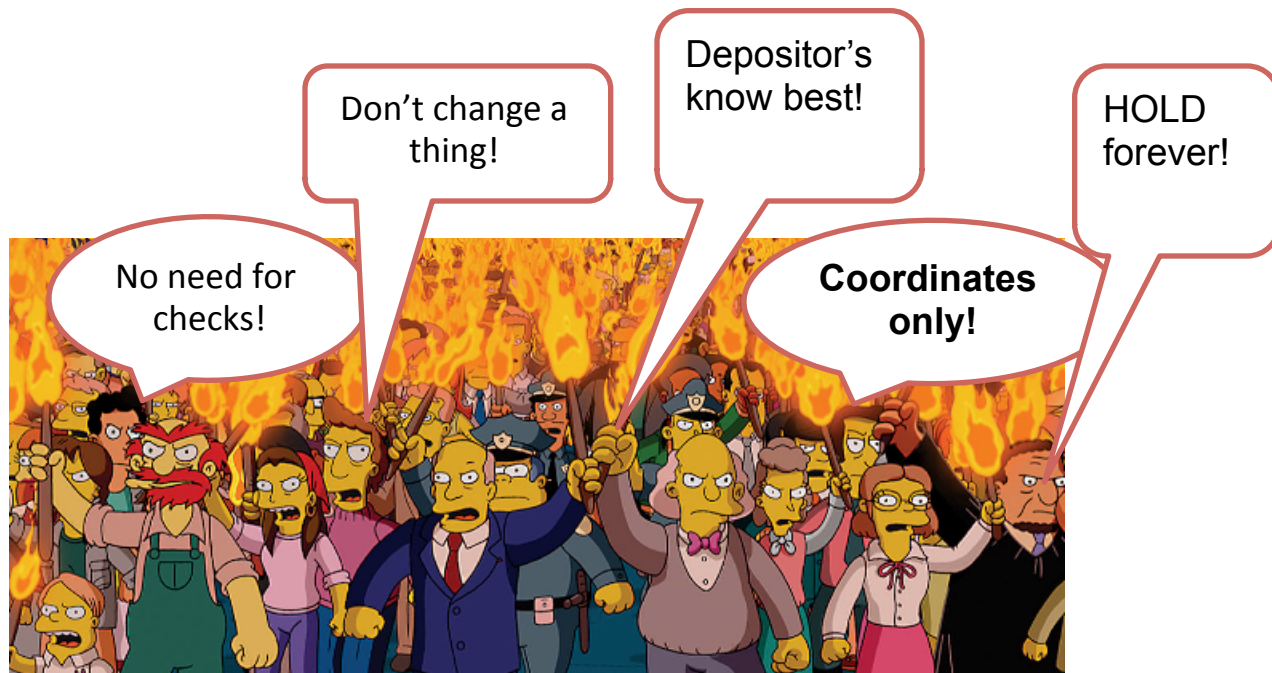- Possible incremental strategy

A unique scientific collaboration providing the authoritative global resource for experimentally determined 3D structures of important macromolecules.

- Formalization of current working practice
- MOU signed July 1, 2003
- Announced in *Nature Structural Biology* November 21, 2003
- **Each partner funded locally**

# Changing View of PDB

## A Generally Hands-Off Role

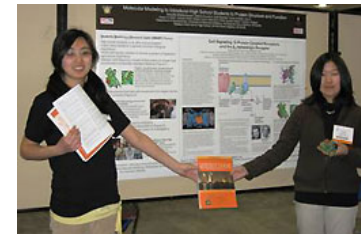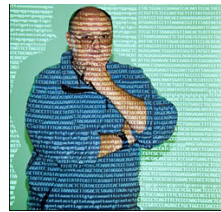# Changing View of  PDB
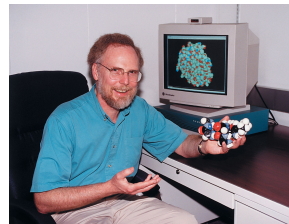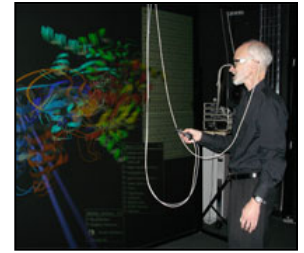
## Increasing Emphasis on Data Quality

# Changing View of PDB

## Increased Emphasis on Data Archiving

# Who Downloads PDB data?

- Structural biologists
- Experimental biologists
- Computational biologists
- Biochemists
- Molecular biologists
- Educators
- Students

# How Do We Interact With These Communities

*Interactions with different user communities*



Feedback

Outreach

biochemists
bioinformaticists teachers
undergraduates
structural biologists
geneticists
high school students
data depositors
software developers
graduate students
illustrators
pharmacologists authors
K12
data users
publishers
textbook

*Development of PDB resources*

# Community Driven Data Standards PDBx/mmCIF

# wwPDB Task Forces

To collect recommendations and develop consensus on method-specific issues, including validation checks that should be performed and identification of validation software applications.

**X-ray Validation**
- 2008 Workshop
- 2011 *Structure* publication
- Chair: Randy J. Read (University of Cambridge)

**3DEM Validation**
- 2010 Meeting
- Chairs: Richard Henderson (Maps, MRC-LMB), Andrej Sali (Models, UCSF)
- White paper in progress

**NMR Validation**
- Meetings held 2009, 2011
- Chairs: Gaetano Montelione (Rutgers), Michael Nilges (Institut Pasteur)
- Report in progress

**Small-Angle Scattering**
- 2012 Meeting
- Members: Jill Trewhella (Univ Sydney), Dmitri Svergun (EMBL Hamburg), Andrej Sali (UCSF), Mamoru Sato (Yokohama City Univ), John Tainer (Scripps)

NMR

X-ray

EM

# Workshops and Working Groups

## 3DEM Data Exchange
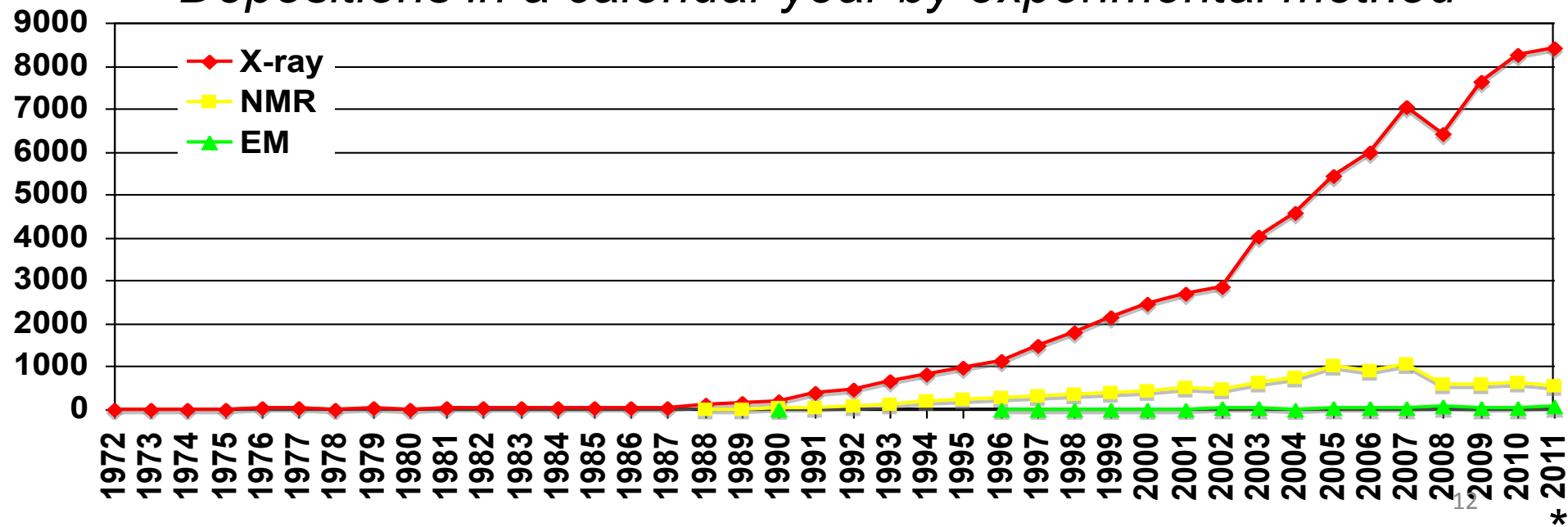I2PC Workshop 2012 - Madrid

## PDBx Deposition Working Group
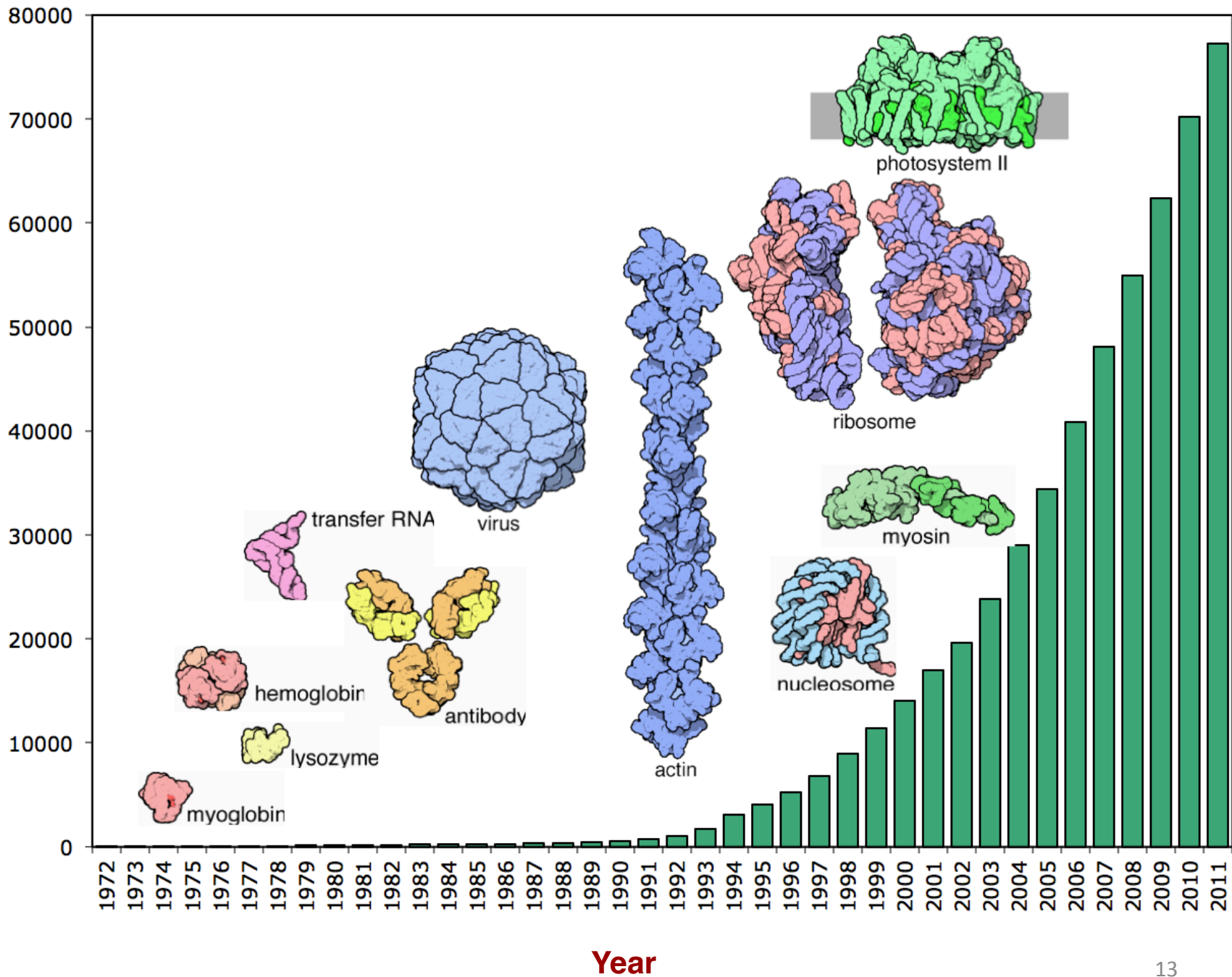Refinement Developers Workshop 2011 - EBI

# 10,000-Fold Growth in Four Decades

- 83,400 entries
- 2012 will see ~10,000 depositions
- Over 85% entries include structure factor data used in the final refinement

*Depositions in a calendar year by experimental method*

The x-axis is labeled "Year" with values from 1972 to 2011. The y-axis is labeled "Number of released entries" ranging from 0 to 80000. The figure includes labeled protein structure images: photosystem II, ribosome, virus, actin, myosin, nucleosome, transfer RNA, hemoglobin, antibody, lysozyme, myoglobin.

# PDB FTP Downloads



Version 3.0/3.1 files released

Version 3.15 files released

Version 4.0 files released

~2 M downloads/year of structure factor data files 2009-present

# 2010 FTP Traffic



| RCSB PDB | PDBe | PDBj |
|----------|------|------|
| 159 million entry downloads | 34 million entry downloads | 16 million entry downloads |

# Challenges and Scope

- Target content
- Longevity
- Example applications
- Audience
- Representation

# What are the Content Targets?

- Laboratory data files
- Laboratory data files with supporting metadata
- Archival storage of standardized data and metadata

# Expected Duration of Storage?

- Through publication review
- A few years not to exceed the availability of supporting hardware & software
- Longer …

# Use Cases

- Recover laboratory data files
- Satisfy philosophical/ethical/funding requirements
- Support peer review, reproduction, and validation of published results
- Extend on published results
- Provide test cases/benchmarks for methods development
- Preserve data from difficult cases

# Audiences  Impacted

- Direct Impact
  - Methods developers
  - Expert users
- Indirect Impact
  - Novice or non-specialist users

# What format and metadata?

# Archival Format and Metadata

- Solid metadata foundation for archiving -
  - CBF/imgCIF
  - PDBx/mmCIF
- Not widely used at early stages of the structure determination pipeline.
- How will working formats and process details be standardized for archiving?
  - Tar ball containing 1000 files from multiple data collections, multiple crystals, with multiple wavelengths …

# Format and Metadata Targets

- Existing efforts provide data in program formats and limited software accessible metadata (e.g. TARDIS & JCSG)

- To what extent does this limit the audience and the useful lifetime of this data?

# Technical Options

- Self-publishing
- Institution/facility hosting and delivery
- Centralized cloud delivery
- Centralized delivery by the PDB

# Self-Publishing

- Contributor posts contents to a file sharing resource
  - Institution or facility storage resources
  - Google Drive –
    - 25GB $2.50/month - 16TB $800/month
  - Egnyte Hybrid Cloud
    - 150GB $300/yr
  - FileSwap
    - Up to 50 GB for $9.95/month
- Contributor registers DOI and digital signatures with archive

# Centralized Cloud Delivery

Target one year ~ 10,000 x 5GB data sets

- ## Leased storage from a major provider
  - Amazon –
    - Storage - $0.125/GB/month
    - Access - ~$0.12-0.05/GB + $0.01/request
  - Google
    - Storage - $0.095/GB/month
    - Access - ~$0.21-0.08/GB + $0.01/request

- ## Application developed to manage depositions

- ## DOIs and signatures registered with archive

$75K storage + $5100/download in yr 1
$450K storage + $15.3K/download after yr 3

# Archive Centralized Storage Hardware Costs

Target one year ~ 10,000 x 5GB data sets

- ## Cheap RAID or JBOD
  - 50TB ~ $30K or ~ $600/TB w/ 3yr maintenance
- ## NAS Expansion (disks and shelves only)
  - NetApp –
    - 50 TB ~ $83K or $1675/TB w/ 3yr maintenance
  - DDN -
    - 50 TB ~ $51K or $1025/TB w/ 3yr maintenance

# Archive Centralized Storage
## Minimum System Requirements

- Deposition site primary and backup copy

- Distribution site primary and backup copy

- Assume data requirement of 50 TB per year for the first 3 years -
  - Cheap RAID - $360 K
  - NetApp expansion - $ 996K
  - DDN expansion - $612 K

- In year 4, replace existing disk hardware + new storage for year four data.

At each wwPDB site

# Archive Curation Costs
## wildly optimistic estimates

- Early Stages –
  - 1 crystallographic application programmer
  - 1-2 annotators with deep expertise and troubleshooting experience with a variety of data collection, integration and phasing applications.
  - 1 scientific programmer to implement deposition and data processing automation

Comparable staffing requirements at each wwPDB site

# Some Possible Practical Steps

- Tackle unmerged intensities first

- Register DOIs and digital signatures for locally store/self-published image data sets.

- Develop metadata extensions for all processing steps.

- Implement standard formats and metadata with facility control systems and pipeline software

- Pilot an automated data capture system with standard data format and metadata.

# Acknowledgements

**Operated by two members of the RCSB:**

**The RCSB PDB is a member of the**

**Supported by:**