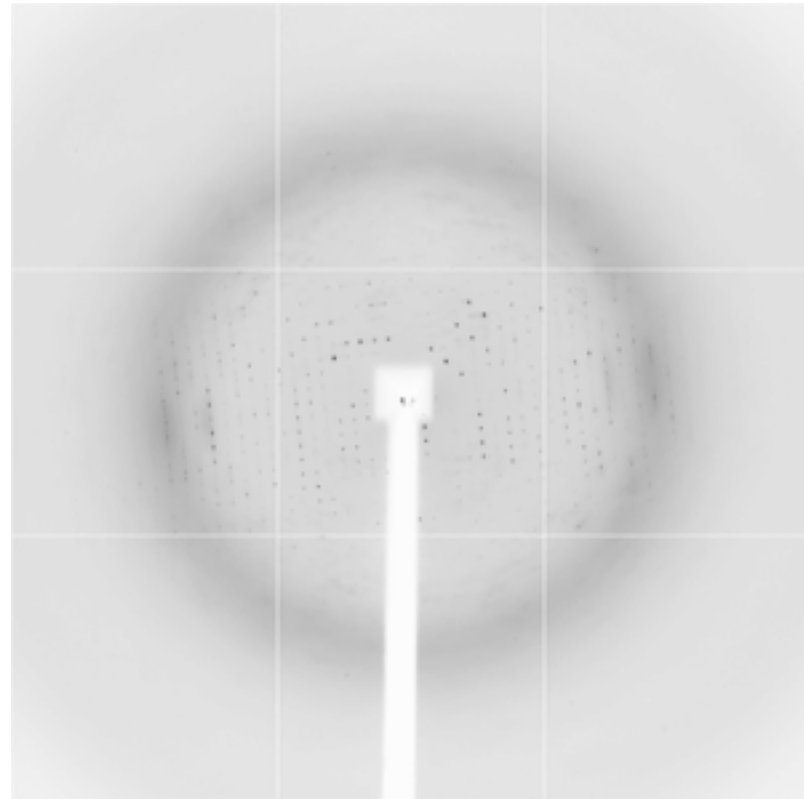National Data Service Pilot

# Establishing an effective system to facilitate access to biomedical datasets.

Peter Meyer
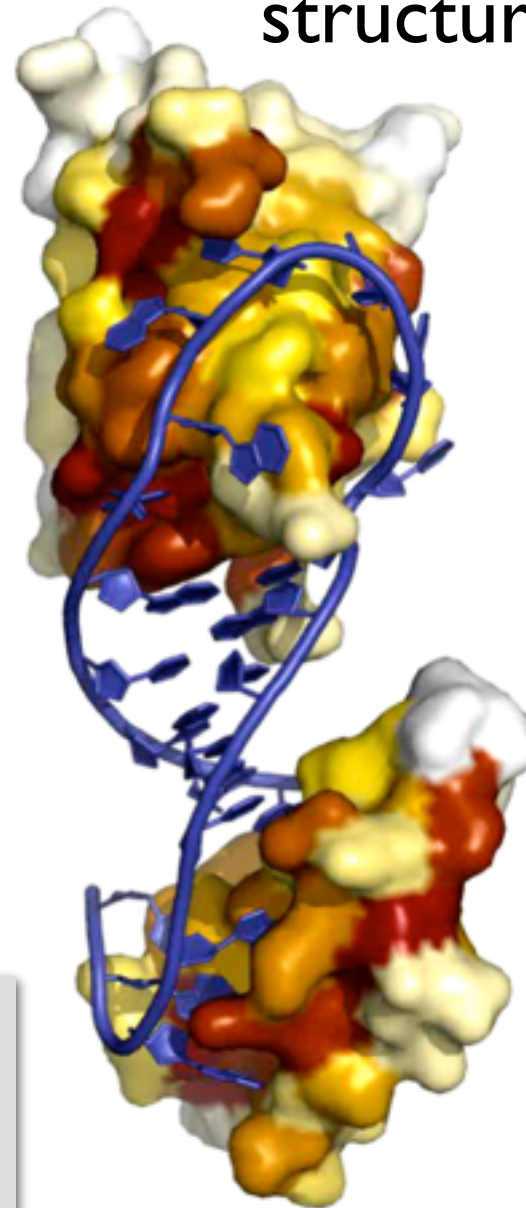
Piotr Sliz, Merce Crosas, Ian Foster

National Data Service
UNC-Chapel Hill, April 5th, 2016
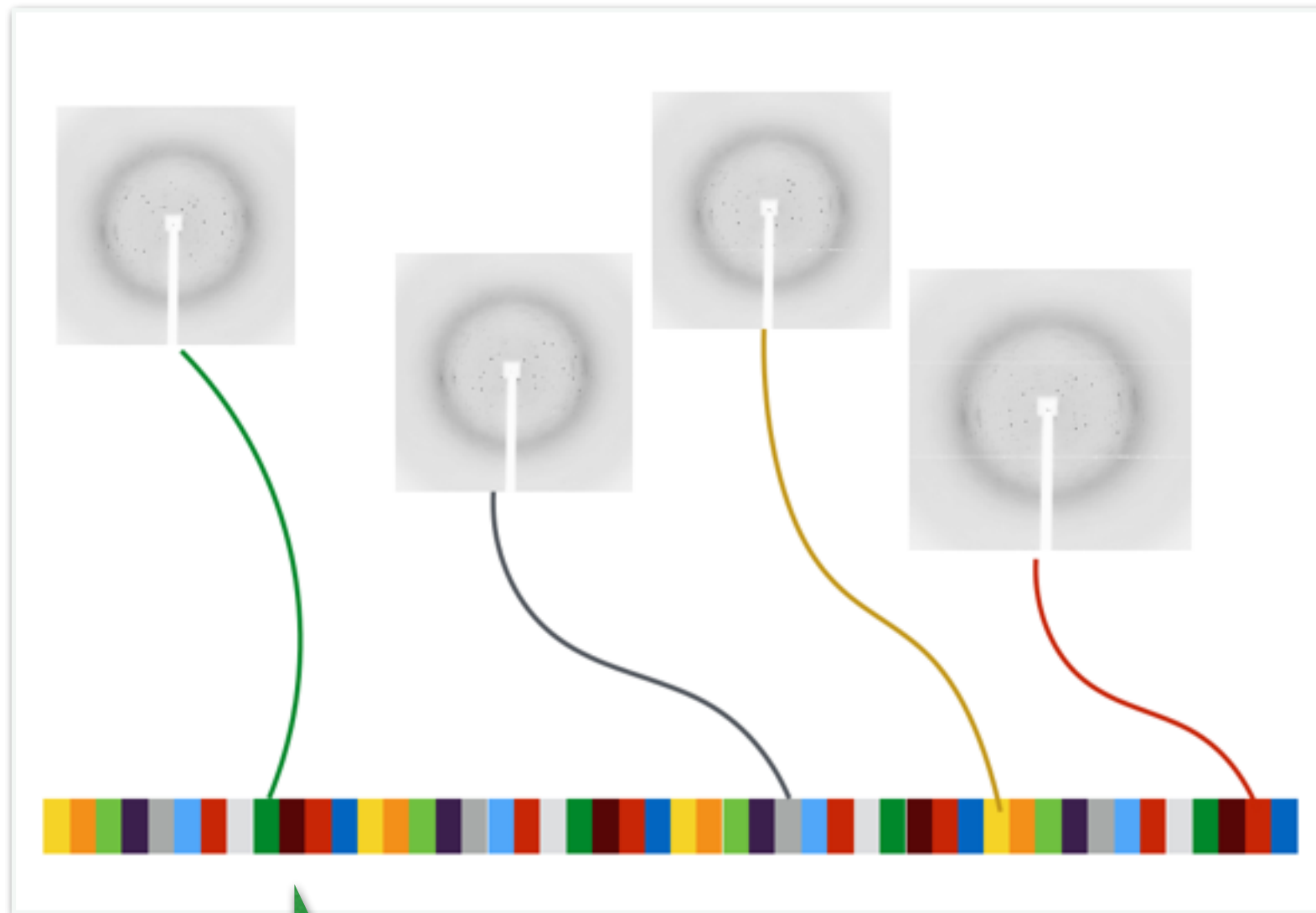
# Motivation



Lin28 structure

MD Models

Nam et al., Cell, 2011

- Different software packages (e.g. XDS vs HKL2000: 3D vs 2D profile fitting)
- Different assumptions (e.g. symmetry, mosaicity, radiation damage, number of frames)
- New software packages (e.g. DIALS)
- Improved criteria (e.g. resolution limits such CC1/2, Karplus and Diederichs, 2012, or anisotropic correction)
- New corrections (e.g. data anisotropicity)
- Additional features: (e.g. anisotropic diffuse scattering signals)

WORLDWIDE
wwPDB
PROTEIN DATA BANK

# Four Essential Components of a biomedical data grid

# Structural Biology Data Grid Website

## data.sbgrid.org

Lab Collections



doi:10.1038/ncomms10882

nature.com    journal home    archive by date    march    abstract

NATURE COMMUNICATIONS | ARTICLE    OPEN

## Data publication with the structural biology data grid supports live analysis

Peter A. Meyer Stephanie Socias Jason Key Elizabeth Ransey Emily C. Tjon Alejandro Buschiazzo Ming Lei Chris Botka James Withrow David Neau Kanagalaghatta Rajashankar Karen S. Anderson Richard H. Baxter Stephen C. Blacklow Titus J. Boggon Alexandre M. J. J. Bonvin Dominika Borek Tom J. Brett Amedeo Caflisch Chung-I Chang Walter J. Chazin Kevin D. Corbett Michael S. Cosgrove Sean Crosson Sirano Dhe-Paganon Enrico Di Cera Catherine L. Drennan Michael J. Eck Brandt F. Eichman Qing R. Fan Adrian R. Ferré-D'Amaré J. Christopher Fromme K. Christopher Garcia Rachelle Gaudet Peng Gong Stephen C. Harrison Ekaterina E. Heldwein Zongchao Jia Robert J. Keenan Andrew C. Kruse Marc Kvansakul Jason S. McLellan Yorgo Modis Yunsun Nam Zbyszek Otwinowski Emil F. Pai Pedro José Barbosa Pereira Carlo Petosa C. S. Raman Tom A. Rapoport Antonina Roll-Mecak Michael K. Rosen Gabby Rudenko Joseph Schlessinger Thomas U. Schwartz Yousif Shamoo Holger Sondermann Yizhi J. Tao Niraj H. Tolia Oleg V. Tsodikov Kenneth D. Westover Hao Wu Ian Foster James S. Fraser Filipe R. N C. Maia Tamir Gonen Tom Kirchhausen Kay Diederichs Mercè Crosas Piotr Sliz

**Affiliations    Contributions    Corresponding author**

*Nature Communications* 7, Article number: 10882    doi:10.1038/ncomms10882
Received 16 October 2015    Accepted 28 January 2016    Published 07 March 2016

| Full text | PDF | Citation | Reprints | Rights & permissions | Article metrics |

### Online attention

81

Altmetric score (what's this?)
Tweeted by 36
Mentioned in 1 Google+ posts
Picked up by 7 news outlets
Blogged by 1

This Altmetric score means that the article is:

• in the 97 percentile (ranked 4,195th) of the 169,738 tracked articles of a similar age in all journals
• in the 80 percentile (ranked 124th) of the 647 tracked articles of a similar age in *Nature Communications*

# Storage Requirements
## (based on 100,000 PDB files and average dataset size in SBDG)



| Per Project: | 100,000 Projects: |
|---|---|
| 0.5 TB | 48 PB |

**All Experimental Data Diffraction Images**

x 6,200,000

6.2 GB → 602 TB

**Primary Data** x 1.26 Dataset/PDB

x 6,200

1 MB → 97 GB

**Structure-factor Amplitudes** x 4

0.23 MB → 22 GB

**Macromolecular Models**

SBDG:
110 datasets
~0.5TB

PDB:
100,000 models
0.3 TB

Some of "All Experimental Data" preserved at national synchrotrons e.g. Tardis or Diamond

Primary Diffraction Datasets proposed to be stored on SB Data Grid

**SBGrid** CONSORTIUM

Molecular models and reduced datasets are stored in Protein Data Bank

WORLDWIDE **PDB** PROTEIN DATA BANK

## Data Access

# Dataset Access Alliance:
## Local access through a growing list of Satellites



UPPSALA UNIVERSITET

Petrel Storage
at Argonne Labs

Orchestra Cluster
Harvard Medical School

San Diego
Supercomputer
Center

INSTITUT PASTEUR
DE MONTEVIDEO

SIBS  SHANGHAI INSTITUTES FOR BIOLOGICAL SCIENCES,
CHNINESE ACADEMY OF SCIENCES

## National Data Service Pilot

Dataset Access Alliance:
Local access through a growing list of Satellites

SBDG Endpoint
Globus Connect Server

transfer

Yale MS Endpoint
Globus Connect Personal

coordination

Control
System

REST API
Calls

Globus
Network

National Data Service Pilot

@SBGrid

# Integration with Biomedical Software Support



sbgrid.org
biogrids.org

**Stephanie Socias**

**Pete Meyer**

# Thank you

*@SBGrid*