SBGrid Databank

Peter A Meyer¹, Stephanie Socias¹, Jason Key¹, Mercè Crosas², Piotr Sliz^{1,3}

BCMP, Harvard Medical School and SBGrid Consortium;
IQSS, Harvard University and Dataverse Project;
Boston Children's Hospital, Dept. of Pediatrics



https://data.sbgrid.org

Problem

- Diffraction images for published structures not widely available
 - Limits re-analysis and validation
 - Difficult to use as controls
 - Limits method development
 - May impact funder data retention policies



Why haven't images been as widely available as PDBs?

- Needs more storage [10s GB vs 10s MB]
- Wider format variations
- Concerns about how reusable images will be
- It's more work for experimenters

What's a dataset?

- Generically
 - Data files: inputs to analysis/compute
 - Metadata: information needed to use data files
- Primary X-ray diffraction datasets
 - Image files
 - Resulting in single Intensity/Amplitude set



Getting Data Out

/1 .

Data Access Instructions	<pre>pameyer@pm-linux:~\$ rsync -av rsync://data.sbgrid.org/10.15785/SBGRID</pre>
1. If this dataset is locally available, it should be accessable at /oroorams/dataorid/1	receiving incremental file list
2. To download this dataset, please run the following command from your Terminal on a Linux or OS X workstation:	1/
'rsync -av rsync://data.sbgrid.org/10.15785/SBGRID/1 .' (Harvard Medical School, USA)	1/files.sha
Depending on your location, faster access may be available from a Tier 1 site closer to your location	1/p3_6_1.runlist
'rsync -av rsync://sbgrid.icm.uu.se/10.15785/SBGRID/1 .' (Uppsala University, Sweden)	1/p3_6_1_001.1mg
'rsync -av rsync://sbgrid.pasteur.edu.uy/10.15785/SBGRID/1 .' (Institut Pasteur de Montevideo, Uruguay)	$1/p_{3} = 1.002.1mg$
'rsync -av rsync://sbgrid.ncpss.org/10.15785/SBGRID/1 .' (Shanghai Institutes for Biological Sciences, China)	1/p3_6_1_004.img
3. After the transfer is completed, please issue the following command to verify data integrity:	1/p3_6_1_005.img
'cd 1 ; shasum -c files.sha'	1/p3_6_1_006.img
Storage requirements: 1.6G	1/p3_6_1_007.1mg
	1/p3_6_1_008.img
	1/p3_6_1_009.1mg
	1/p3_6_1_010.1mg
	$1/p_{2} = 1.012$ imp
	$1/p_{3} = 1.012$, img
	$1/p_{3} = 0.12$, lmg
	1/p3_0_1_01+, ting
	pameyer@pm-linux:/programs/datagrid/1\$ shasum -c files.sha
	./p3_6_1_006.img: OK
	./p3_6_1_001.img: OK
	./p3_6_1_073.img: OK
	./p3_6_1_008.img: OK
	./p3_6_1_074.img: OK
	./p3_6_1_090.img: OK
	./p3_6_1_041.img: OK
	./p3_6_1_046.img: OK
	./p3_6_1_048.img: OK
	./p3_6_1_034.img: OK
	./p3_6_1_033.img: OK
	./p3_6_1_022.img: OK
	./p3_6_1_059.img: OK
	./p3_6_1_088.img: OK
	./p3_6_1_025.img: OK

Putting Data In

Data 🔻	About 🔻	Get Help 🔻	For Depositors v	Q
			Request Deposition Credentials	
			Register Dataset	

Depositor Name: **Depositor Email: Depositor Institution:** PI Name: PI Email: PI Institution: PI ORCID: Lab website: This is a request from a new lab or group: Captcha:

• Date Collected/Cre mm/dd/yyyy	ated*	Subject Composition*
Date Collected/Cre mm/dd/yyyy	ated*	Subject Composition*
mm/dd/yyyy	86	C PMI C Useral C Particip
		DNA Ligand Protein RNA
		0 max
	Dataset loop	
	Browse No fi	le selected.
	Publication DOI	
ndicata if sha/ha als	o acted as a data co	lactor or depositor
receive the email w	ith dataset upload in	istructions.
Middle initial		Last name*
	Email*	
	ndcate if she/he also receive the ernal is Niddle initial	Dataset Icon Browse No fi Publication DOI Indicate if sho/he also acted as a data oc receive the email with dataset upload in Widdle Initial

UPLOAD DATA

Register your Data

Agree to the Terms and Conditions

Cancel

1

REGISTER DATA

eyer@pm-linux:~\$ bash ./upload2sbdb_dataset-No0.bash

ease enter the full path of the directory containing dataset files to upload for dataset number: 0 This directory should contain only the dataset in question - all files within this directory will be uploaded /data/images/

/data/images/p3_6_1_001.img /data/images/p3_6_1_002.img data/images/p3_6_1_003.img data/images/p3_6_1_004.img /data/images/p3_6_1_005.img /data/images/p3_6_1_006.img data/images/p3_6_1_007.img data/images/p3_6_1_008.img data/images/p3_6_1_009.img

ta/images/p3_6_1_010.img

1

PUBLICATION COMPLETE

Request Credentials

X-Ray Diffraction Collection Changes

	NC 20150901	20170524
Datasets	111	311
PDBs	91	274
Articles	65	128
Labs / Groups	48	66
Institutions	34	42
Collection Facilities*	65	151
Detectors*	8	23

How reusable is this data?

- Crystallographic structure determination is multi-step
- Reusability needs to evaluated stepwise
- With known structures, can be evaluated "backwards"





309 datasets

How to get better metadata?

- Encourage and support standardized formats with accurate experimental metadata
- Ask depositors to provide experimental parameters
- Have curators attempt to to determine experimental metadata
- Have computers attempt to determine experimental metadata

What to do with new metadata?

- "Correct" images
 - Complicates data integrity
 - Headers not designed for multiple sources of information
- REST API
 - Complicates analysis pipelines







control RFree

image RFree

reprocessed

Connecting to Scientific Ecosystem

Cite this Dataset

Nam, Y; Sliz, P. 2015. "X-Ray Diffraction data for: Lin28A/let-7g microRNA complex. PDB Code 3TS2", SBGrid Data Bank, V1, http://dx.doi.org/10.15785/SBGRID/1.



Connecting to Scientific Ecosystem

- Integrations
 - RCSB integration
 - DataMed
- API



Q

https://data.sbgrid.org/static/html/apidoc/index.html

What's next?

- Dataverse integration
 - Improved UI
 - ORCID login



- Improved analysis pipelines
 - Additional tools
 - API generalization
- User suggestions

Acknowledgements

Stephanie Socias Jason Key Bill McKinney Mick Timony Justin O'Connor Michelle Ottaviano Piotr Sliz

Phil Durbin Gustavo Durand Leonid Andreev Mercè Crosas

Developers / Projects:

- xia2
- CCP4
- XDS
- MOSFLM
- PHENIX
- LABELIT

Crystallographers

Alejandro Buschiazzo Ming Lei Filipe Maia



Protein Data Bank

THE LEONA M. AND HARRY B. HELMSLEY CHARITABLE TRUST