



Crystallographic raw data: our plans and implementations within the NIH's Big Data to Knowledge resource

Wladek Minor

Rovinj, August 2015

Integrated Resource for Reproducibility in Macromolecular Crystallography

Goals:

1. Develop tools for automatically extracting and curating diffraction images and associated metadata, as well as producing detailed descriptions of all data needed for later reprocessing of the diffraction data as methods for structure determination improve.
2. Create a web-based system for semantic searching, analysis, and data mining of appropriate subsets of diffraction images and associated metadata.



- 3. Develop tools to automatically validate, preprocess, and score diffraction images, and to detect potential issues and errors.
- 4. Creation of a repository for diffraction data that did not yield an X-ray structure with the currently available methods.
- 5. Set up a pilot resource incorporating the tools developed in Aims 1-4 to collect a test set of data for development of tools and algorithms for validation and error detection.



Integrated Resource for Reproducibility in Macromolecular Crystallography

Goal:

Protein Crystallography (Structural Biology)
with Speed and Finesse

100 years and still going strong



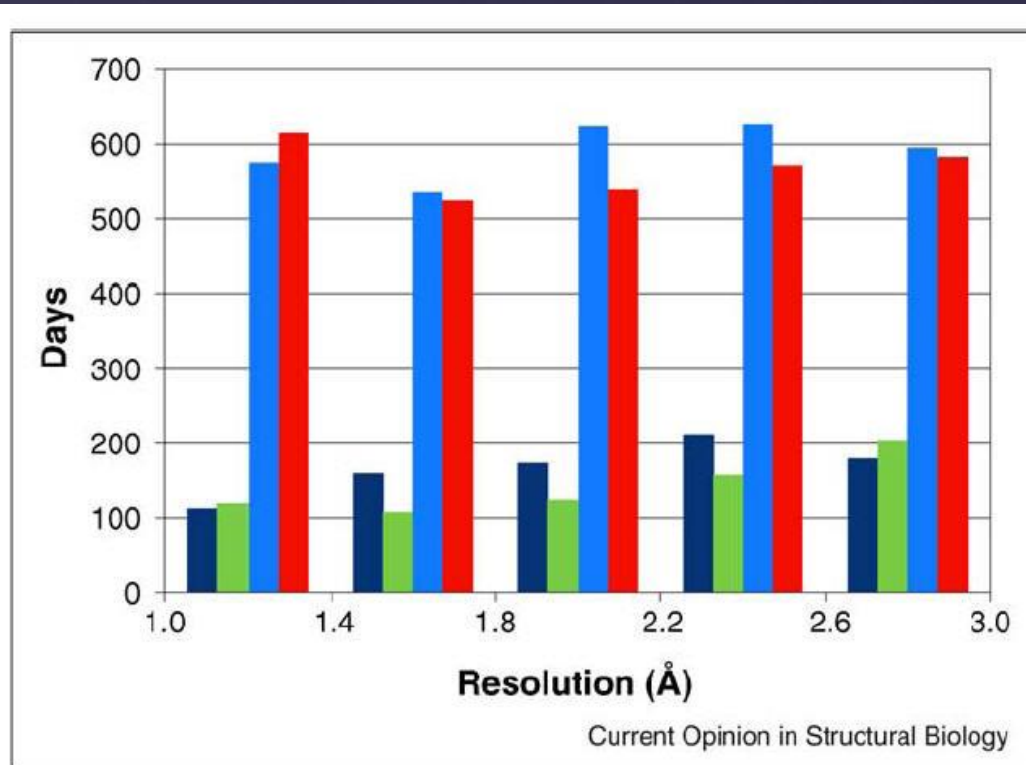
High throughput SB

- Automatic cloning
- HT automatic expression
- HT automatic purification
- HT automatic crystallization
- HT automatic data collection
- HT automatic structure solution/refinement

High throughput SB

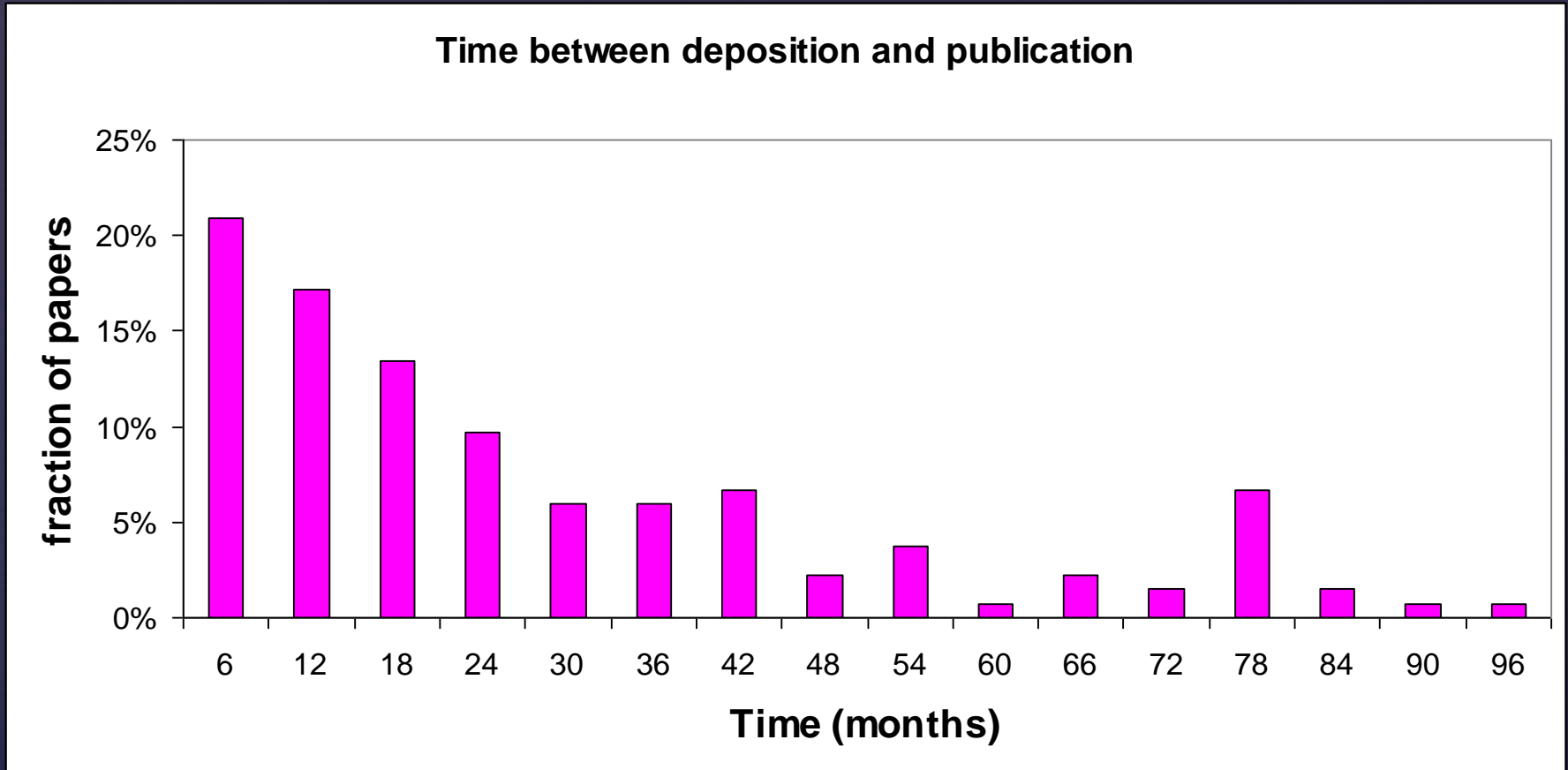
- Automatic cloning
- HT automatic expression
- HT automatic purification
- HT automatic crystallization
- HT automatic data collection
- HT automatic structure solution/refinement
- Automatic paper writing

Data collection -> deposition



Average time (in days) between data collection and deposition for SG and non-SG structures. Dark blue and green bars represent SG structures, whereas light blue and red bars represent non-SG structures deposited in 2000–2004 and 2005–2009, respectively. Structures were binned by reported resolution limit (0.4 Å bin width).

Deposition -> Publication



Bottleneck:

Data is not information,

information is not knowledge,

knowledge is not understanding,

understanding is not wisdom

Bottleneck: brain engagement

**Data is not information,
information is not knowledge,
knowledge is not understanding,
understanding is not wisdom**

Clifford Stoll

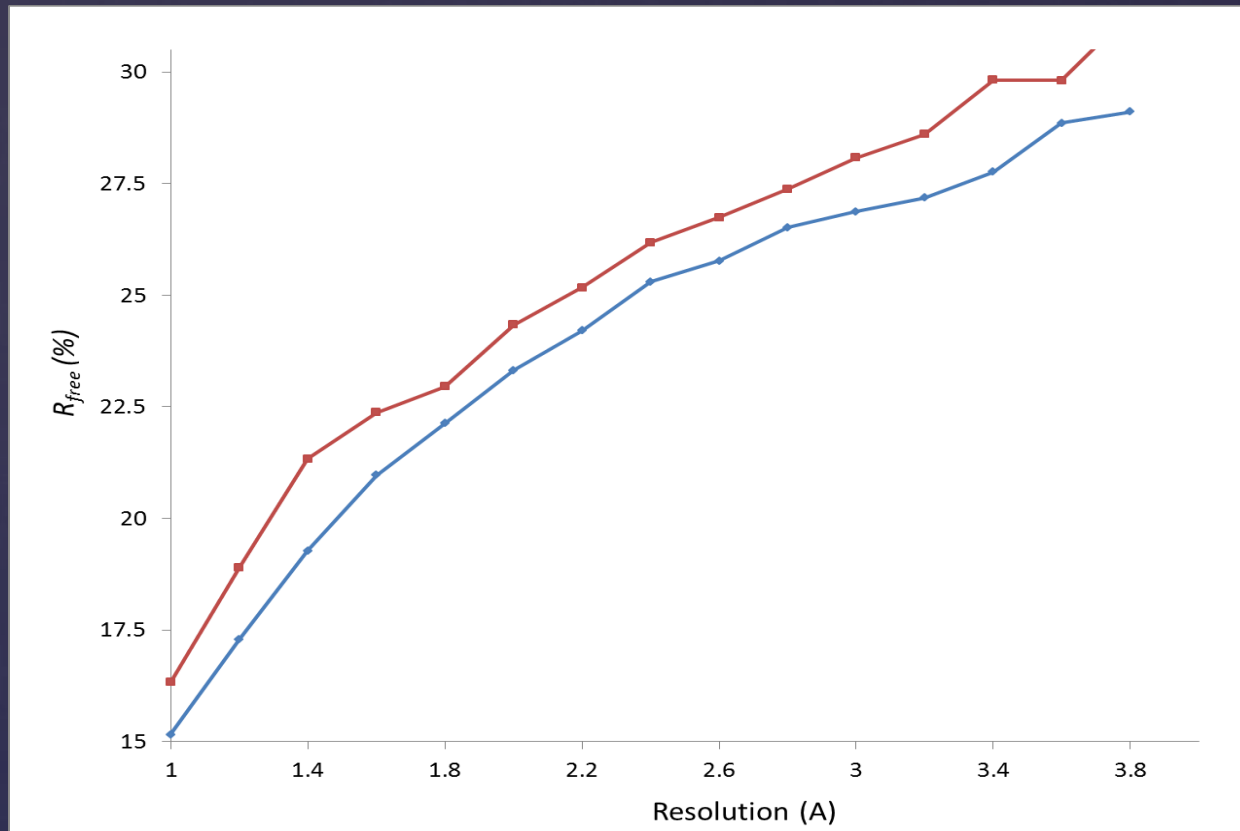
What experimenters know about data collection ?

```
REMARK 3 ESTIMATED OVERALL COORDINATE ERROR.
REMARK 3 ESU BASED ON R VALUE (A) : NULL
REMARK 3 ESU BASED ON FREE R VALUE (A) : NULL
REMARK 3 ESU BASED ON MAXIMUM LIKELIHOOD (A) : NULL
REMARK 3 ESU FOR B VALUES BASED ON MAXIMUM LIKELIHOOD (A**2) : NULL
REMARK 3
REMARK 3 RMS DEVIATIONS FROM IDEAL VALUES.
REMARK 3 DISTANCE RESTRAINTS. RMS SIGMA
REMARK 3 BOND LENGTH (A) : NULL ; NULL
REMARK 3 ANGLE DISTANCE (A) : NULL ; NULL
REMARK 3 INTRAPLANAR 1-4 DISTANCE (A) : NULL ; NULL
REMARK 3 H-BOND OR METAL COORDINATION (A) : NULL ; NULL
REMARK 3
REMARK 3 PLANE RESTRAINT (A) : NULL ; NULL
REMARK 3 CHIRAL-CENTER RESTRAINT (A**3) : NULL ; NULL
REMARK 3
REMARK 3 NON-BONDED CONTACT RESTRAINTS.
REMARK 3 SINGLE TORSION (A) : NULL ; NULL
REMARK 3 MULTIPLE TORSION (A) : NULL ; NULL
REMARK 3 H-BOND (X...Y) (A) : NULL ; NULL
REMARK 3 H-BOND (X-H...Y) (A) : NULL ; NULL
REMARK 3
REMARK 3 CONFORMATIONAL TORSION ANGLE RESTRAINTS.
REMARK 3 SPECIFIED (DEGREES) : NULL ; NULL
REMARK 3 PLANAR (DEGREES) : NULL ; NULL
REMARK 3 STAGGERED (DEGREES) : NULL ; NULL
REMARK 3 TRANSVERSE (DEGREES) : NULL ; NULL
REMARK 3
REMARK 3 ISOTROPIC THERMAL FACTOR RESTRAINTS. RMS SIGMA
REMARK 3 MAIN-CHAIN BOND (A**2) : NULL ; NULL
REMARK 3 MAIN-CHAIN ANGLE (A**2) : NULL ; NULL
REMARK 3 SIDE-CHAIN BOND (A**2) : NULL ; NULL
```

What experimenters know about data collection ?

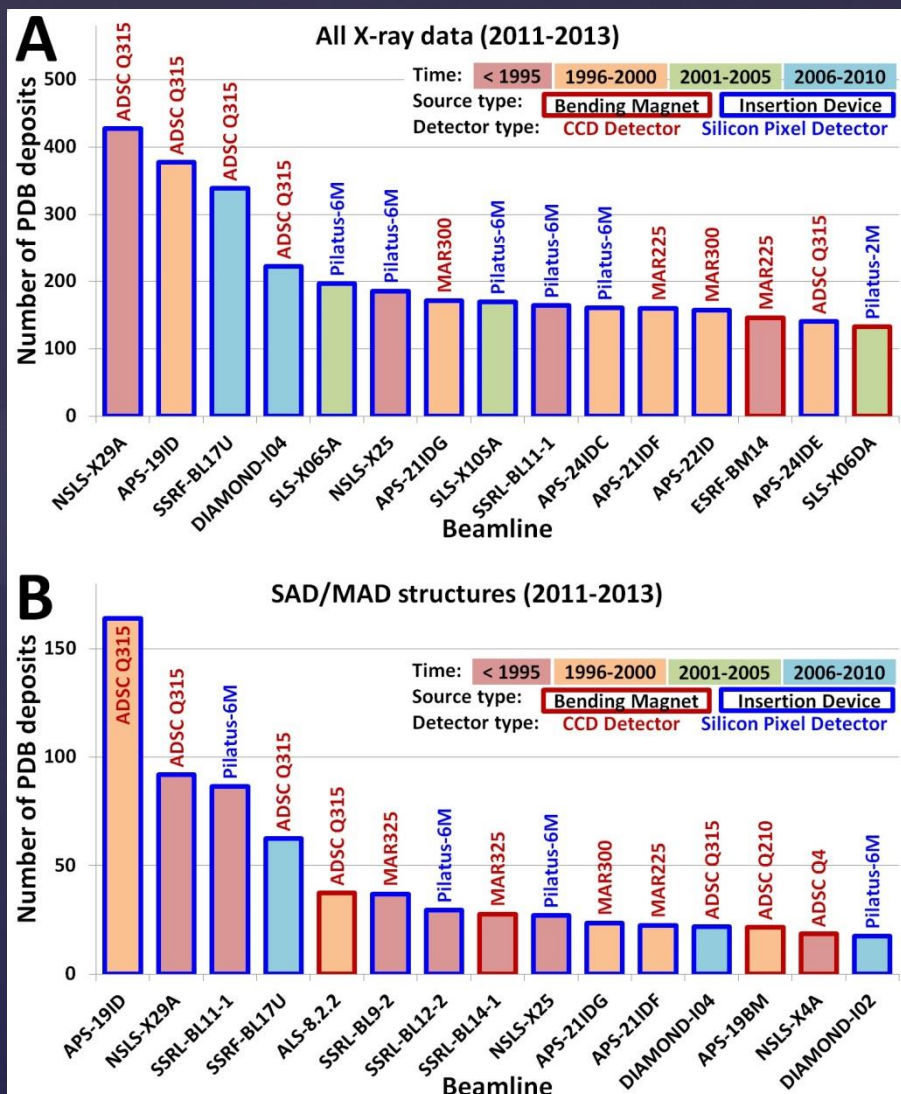
```
REMARK 200 DETECTOR TYPE : CCD
REMARK 200 DETECTOR MANUFACTURER : ADSC QUANTUM 4
REMARK 200 INTENSITY-INTEGRATION SOFTWARE : BLU-ICE
REMARK 200 DATA SCALING SOFTWARE : MOSFLM, CCP4, SCALEPACK
REMARK 200
REMARK 200 NUMBER OF UNIQUE REFLECTIONS : 17575
REMARK 200 RESOLUTION RANGE HIGH (Å) : 2.950
REMARK 200 RESOLUTION RANGE LOW (Å) : 47.870
REMARK 200 REJECTION CRITERIA (SIGMA(I)) : NULL
REMARK 200
REMARK 200 OVERALL.
REMARK 200 COMPLETENESS FOR RANGE (%) : 94.3
REMARK 200 DATA REDUNDANCY : 3.200
REMARK 200 R MERGE (I) : 0.07800
REMARK 200 R SYM (I) : 0.08600
REMARK 200 <I/SIGMA(I)> FOR THE DATA SET : NULL
REMARK 200
REMARK 200 IN THE HIGHEST RESOLUTION SHELL.
REMARK 200 HIGHEST RESOLUTION SHELL, RANGE HIGH (Å) : 2.95
REMARK 200 HIGHEST RESOLUTION SHELL, RANGE LOW (Å) : 3.03
REMARK 200 COMPLETENESS FOR SHELL (%) : 92.4
REMARK 200 DATA REDUNDANCY IN SHELL : 2.80
REMARK 200 R MERGE FOR SHELL (I) : NULL
REMARK 200 R SYM FOR SHELL (I) : NULL
REMARK 200 <I/SIGMA(I)> FOR SHELL : NULL
REMARK 200
REMARK 200 DIFFRACTION PROTOCOL: SINGLE WAVELENGTH
REMARK 200 METHOD USED TO DETERMINE THE STRUCTURE: MAD SE-MET
REMARK 200 SOFTWARE USED: SNB, MLPHARE, CCP4, SOLVE, HKL, RESOLE
REMARK 200 STARTING MODEL: NULL
REMARK 200
REMARK 200 REMARK: NULL
REMARK 280
REMARK 280 CRYSTAL
REMARK 280 SOLVENT CONTENT, VS (%) : NULL
REMARK 280 MATTHEWS COEFFICIENT, VM (ANGSTROMS**3/DA) : NULL
REMARK 280
```


Unexpected correlation?



Average R_{free} by resolution bin (with a width of 0.2 Å for X-ray crystallography PDB structures deposited after January 1, 2001, divided into two groups by the number of missing data items (“NULLs”) in the PDB file. The means for “high-completion” deposits (20 NULLs or less) are shown in blue, and the means for “low-completion” deposits (50 or more NULLs) are shown in red.

Where we should collect data ?



Diffraction experiment - the last experiment before deposition to PDB

Dataset – 2minutes, sample change 2minutes -> 10minutes

6 datasets/hour -> 144 datasets/day

180 days -> 25920 datasets/day -> **2.5 PDB**

125 synchrotron stations -> **324 PDB**

Efficiency -> **0.3%**

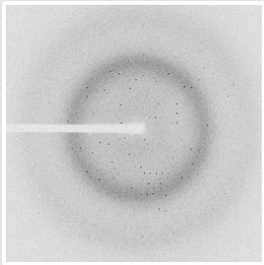
http://www.proteindiffraction.org

http://www.proteindiffraction.org/ x 13 Google Calendar - Week of... x myAT&T Login - Pay Bills ... x 17 Google Calendar - Month ... x +

www.proteindiffraction.org Search

Most Visited http://maps.google.co... Locate Service Center Latest Headlines http://www.lufthansa... Received Messages | Li... Save to Mendeley

Home About Browse Statistics Submit data



Integrated Resource for Reproducibility in Macromolecular Crystallography

This project is being funded by the [Targeted Software Development](#) award 1 U01 HG008424-01 as part of the [BD2K \(Big Data to Knowledge\)](#) program of the National Institute of Health. The project is developing tools for "wrangling" protein diffraction data. We are also creating a growing repository of diffraction images used to determine protein structures in the [PDB](#), contributed by the [CSGID](#), [SSGCID](#), [JCSG](#), [MCSG](#), and other large-scale projects, as well as individual research laboratories.

Currently indexed datasets: 2719

[Read more...](#)

Search examples

Find a specific PDB ID: [4K6A](#)


Free format search: [potential drug target](#)


Combining searches: [drug AND cholera](#)


Specific beamline: [beamline=21-ID-G](#)

Resolution limit: [resolution<1.25](#)

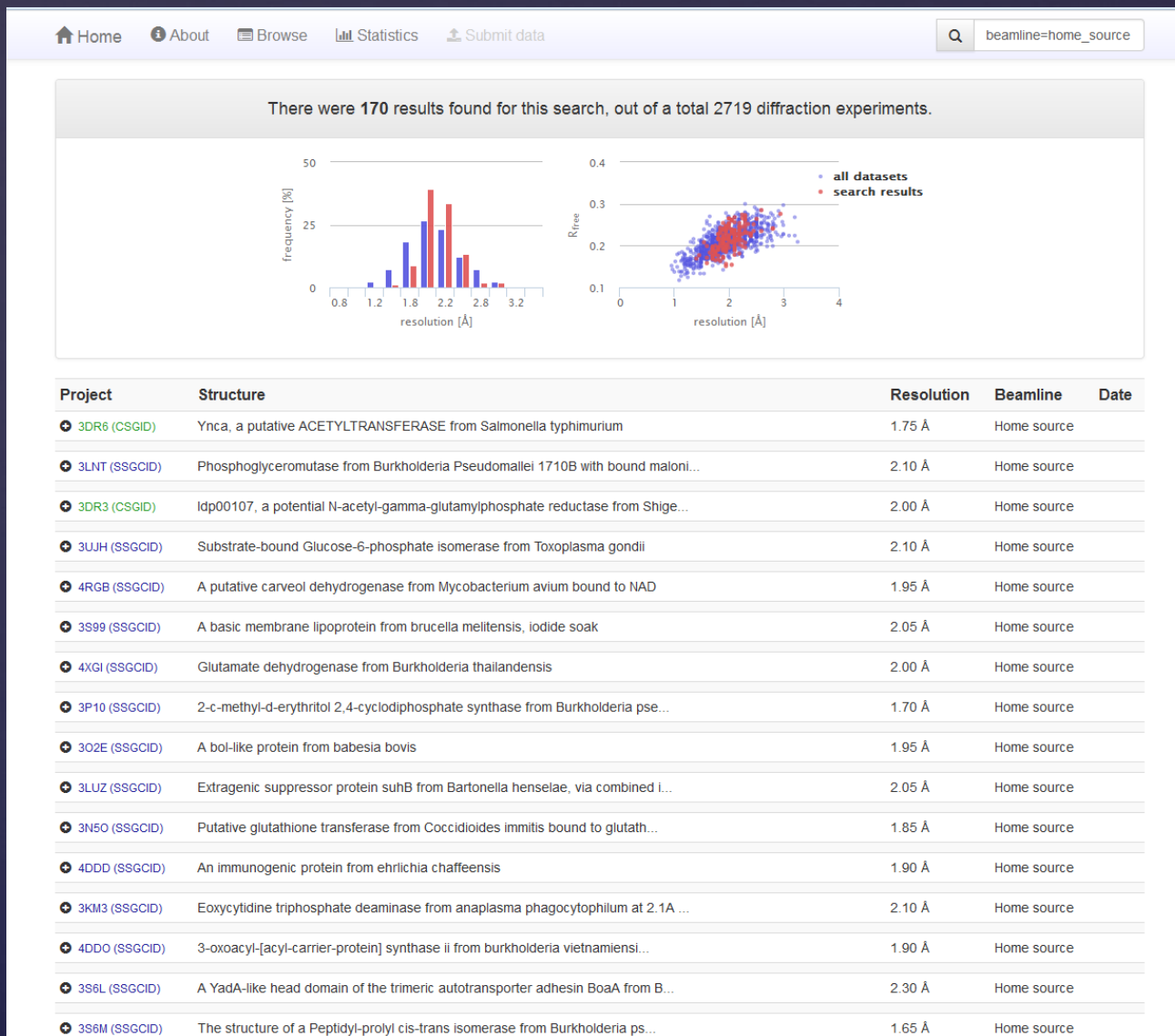
Search by tag: [workshop](#)


Browse & search

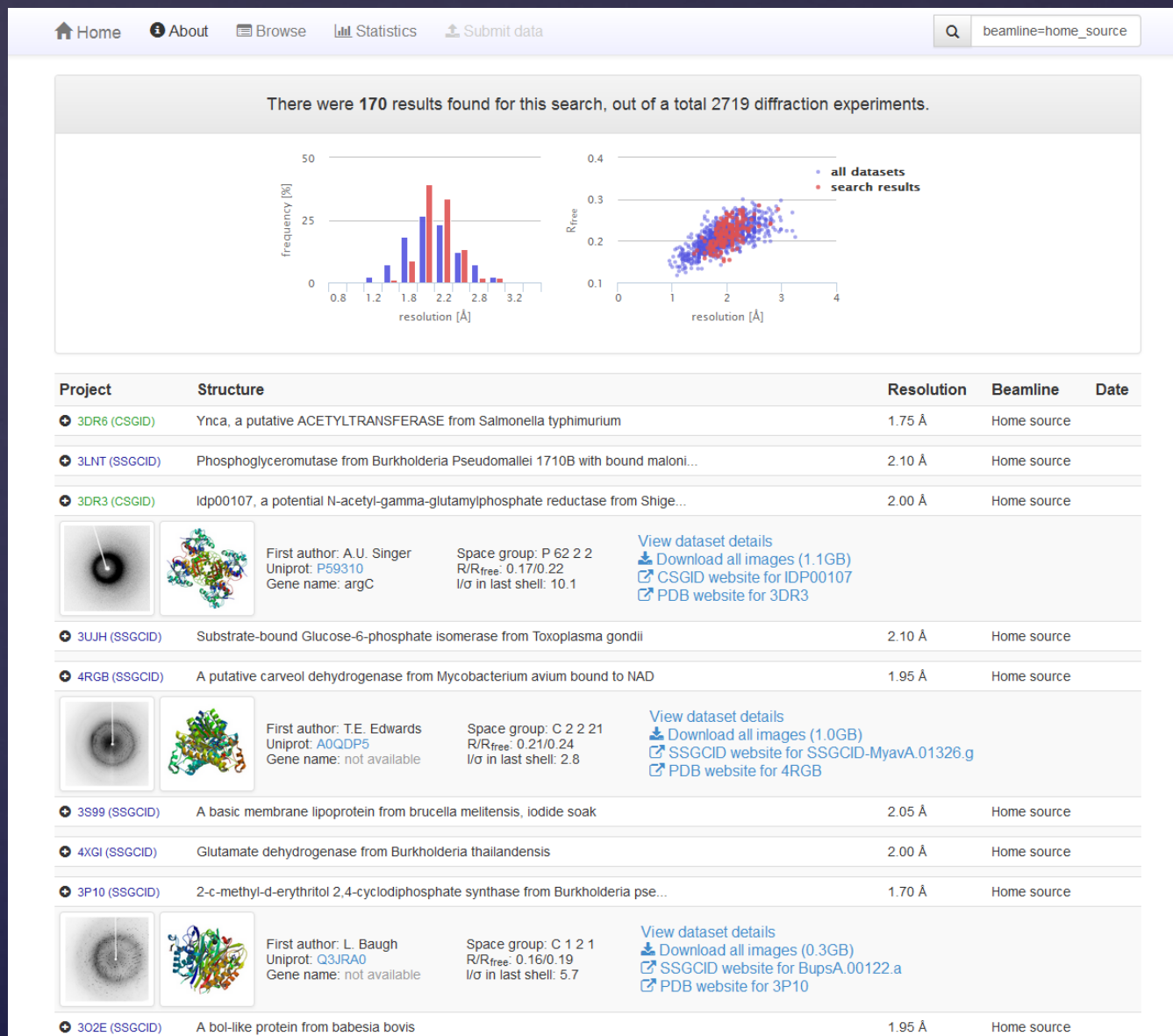

Statistics


Submit data

http://www.proteindiffraction.org



http://www.proteindiffraction.org



Target status and path to success

Space Tree - Mozilla Firefox

File Edit View History Bookmarks Tools Help

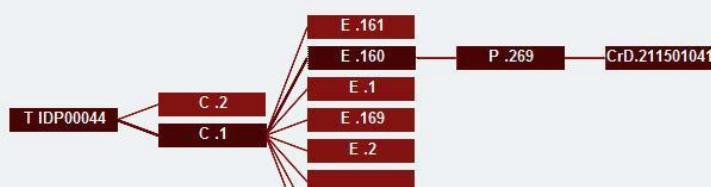
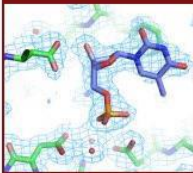
http://csgid.org/csgid/cake/space_tree/view/IDP00044

Most Visited Getting Started Latest Headlines http://www.lufthansa...

Center for Structural Genomics of Infectious Diseases

Home | Target List | Selection | Community Requests | XML files | Diffraction Images | Progress | Homolog Search | Statistics | Help

- Consortium
- Overview
- Investigators
- Targets
- 3-D Structures
- Publications
- Related Sites
- Progress
- Collaborators
- Basecamp
- Log In



Expression
(simple view)

Target	IDP00044
Expression Suffix	160
Protein Clone	1434
Local Clone ID	9
Local Protein Target ID	595
Local Expression ID	160
Experiment Date Start	2008-09-29
Experiment Date End	2008-09-30
Person	Majka Klimecka
Lab	University of Virginia
Protocol ID	uva_expression_native_1
Status	success
Expression Organism	E. coli
Expression Strain	BL21 Codon Plus RILT
Media Type	LB
Media Volume	4 mL
Growth Temperature	37 °C
Induction Temperature	20 °C
Induction Reagent	0.5 mM IPTG
Experiment Type	production
Expression Level	80 %
Solubility Level	80 %

For any help please contact support@csgid.org.

CAKEPHP POWER

Metal binding site validation: *CheckMyMetal* server

ID	Res	Atom	Valence	BV symmetry	Geometry	RMSD geometry angles	Missing vertices	Bidentate	CBVS	Alternative metal
400:A	_MG	MG--	2.1	0.111	Octahedral	5.36	0	0	4.41	
400:B	_MG	MG--	2.06	0.079	Octahedral	4.13	0	0	4.32	

Mouse click action:

☒ None
 ☐ Center
 ☐ Distance
 ☐ Label

Basic controls:

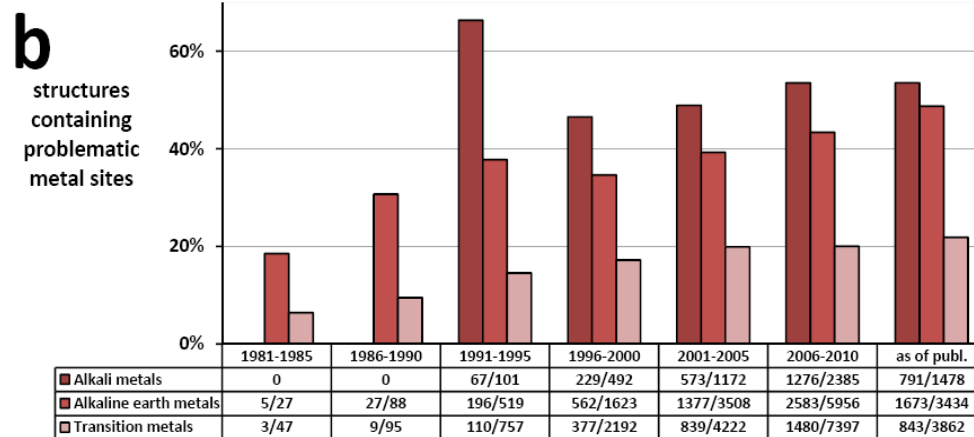
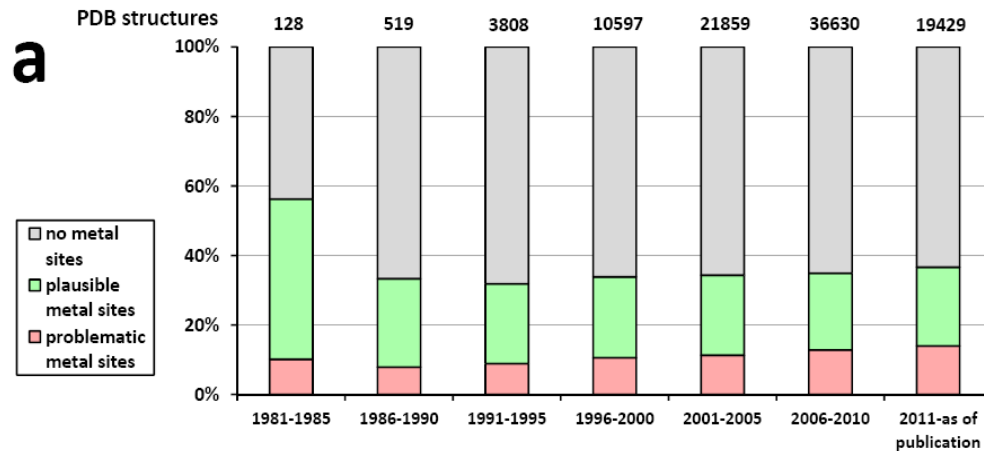
Left-Click to rotate
 Shift-Left-Click up & down to zoom
 Right-Click for Jmol's context menu

Use the buttons below to control the view

Residue Name: ☒ On ☐ Off
 Metal Distances: ☐ On ☒ Off
 Protein Cartoon: ☐ On ☒ Off
 Spin: ☐ On ☒ Off
 Antialiasing: ☒ On ☐ Off

Legend	Explanation
Valence	Summation of bond valence values for an ion binding site
BV symmetry	Summation of bond valence vectors, weighted by bond valence values. Increase when the coordination sphere is not symmetrical due to incompleteness.
Geometry	Ion binding site geometry, as calculated by the NEIGHBORHOOD algorithm
RMSD geometry angles	R.M.S. Deviation of observed geometry angles (L-M-L angles) compared to ideal geometry, in degrees
Missing vertices	Number of sites that is not occupied in the coordination sphere for the assigned geometry
Bidentate	Number of residues that form a bidentate interaction instead of being considered as multiple ligands
CBVS	Calcium Bond Valence Sum, used for alternative metal(s) prediction [Muller, P. et al. (2003) Is the bond-valence method able to identify metal atoms in protein structures? Acta Crystallogr. D Biol. Crystallogr., 59, 32-37.]
Alternative metal	A list of alternative metal(s) is proposed in descending order of confidency, assuming metal environment is accurately determined

Metals in PDB

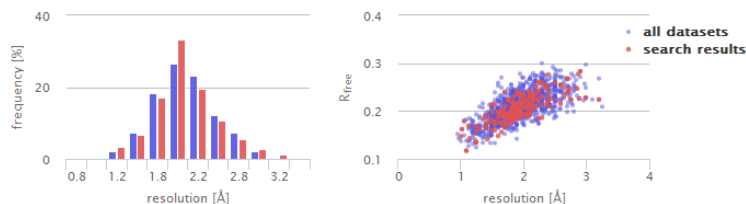


http://www.proteindiffraction.org

Home About Browse Statistics Submit data

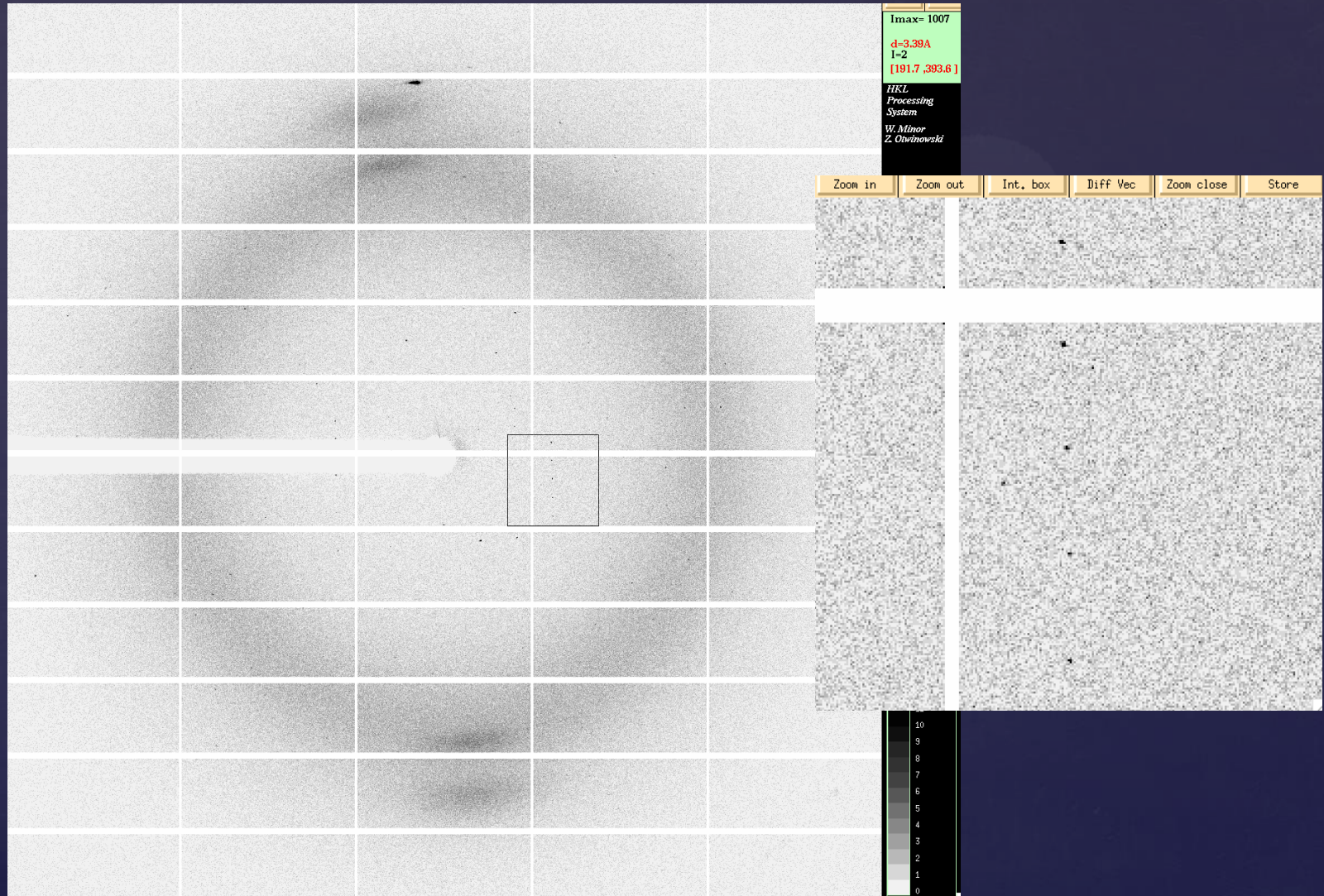
Q beamline=21-ID-G

There were 180 results found for this search, out of a total 2719 diffraction experiments.

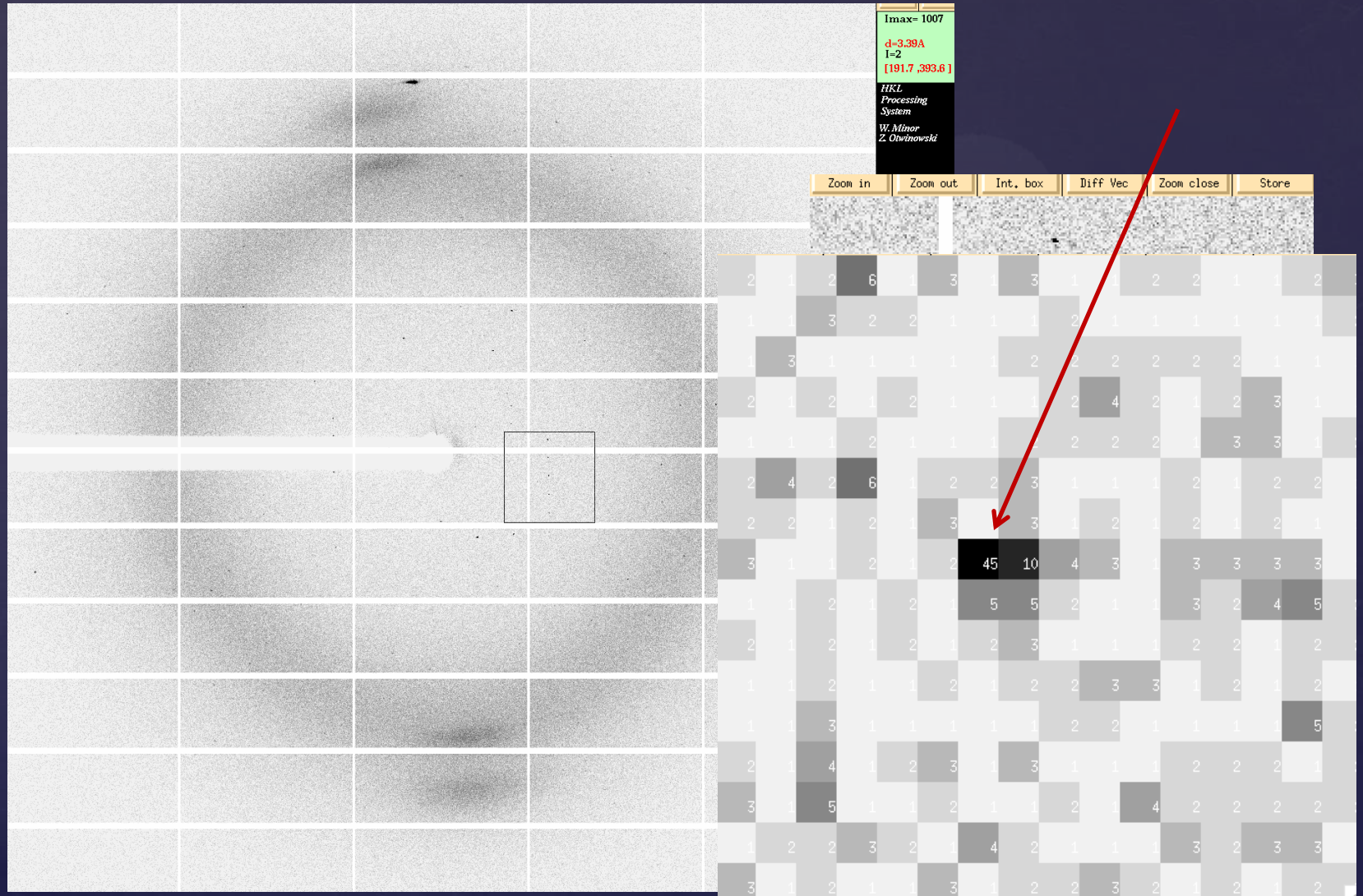


Project	Structure	Resolution	Beamline	Date
4EQ9 (CSGID)	1.4 Angstrom Crystal Structure of ABC Transporter Glutathione-Binding Prote...	1.40 Å	APS / 21-ID-G	
4X9K (CSGID)	Beta-ketoacyl-acyl carrier protein synthase III-2 (FabH2)(C113A) from Vibri...	1.61 Å	APS / 21-ID-G	
4EAQ (MCSG)	Thymidylate Kinase from Staphylococcus aureus in complex with 3'-Azido-3'-D...	1.85 Å	APS / 21-ID-G	
3QTB	Universal stress protein from Archaeoglobus fulgidus in complex with dAMP	2.10 Å	APS / 21-ID-G	
4JBE (MCSG)	1.95 Angstrom Crystal Structure of Gamma-glutamyl phosphate Reductase from ...	1.95 Å	APS / 21-ID-G	
4KWT (CSGID)	Unliganded anabolic ornithine carbamoyltransferase from Vibrio vulnificus a...	1.86 Å	APS / 21-ID-G	
4GIB (CSGID)	2.27 Angstrom Crystal Structure of beta-Phosphoglucosyltransferase (pgmB) from Cios...	2.27 Å	APS / 21-ID-G	
4JG9 (CSGID)	X-ray Crystal Structure of a Putative Lipoprotein from Bacillus anthracis	2.42 Å	APS / 21-ID-G	
4OC9 (CSGID)	2.35 Angstrom resolution crystal structure of putative O-acetylhomoserine (...)	2.35 Å	APS / 21-ID-G	
4HVN (MCSG)	Hypothetical protein with ketosteroid isomerase-like protein fold from Cate...	1.95 Å	APS / 21-ID-G	
3NNT (CSGID)	K170m Mutant of Type I 3-Dehydroquinate Dehydratase (aroD) from Salmonella ...	1.60 Å	APS / 21-ID-G	
3M07 (CSGID)	1.4 Angstrom Resolution Crystal Structure of Putative alpha Amylase from Sa...	1.40 Å	APS / 21-ID-G	
3LAY (CSGID)	Alpha-Helical barrel formed by the decamer of the zinc resistance-associate...	2.70 Å	APS / 21-ID-G	
4OJ7 (SSGICD)	Chorismate Mutase from Burkholderia thailandensis	2.15 Å	APS / 21-ID-G	
3TMQ (SSGICD)	A 2-dehydro-3-deoxyphosphooctonate aldolase from Burkholderia pseudomallei ...	2.10 Å	APS / 21-ID-G	
3IJ3 (CSGID)	1.8 Angstrom Resolution Crystal Structure of Cytosol Aminopeptidase from Co...	1.80 Å	APS / 21-ID-G	

Optimal data collection ?

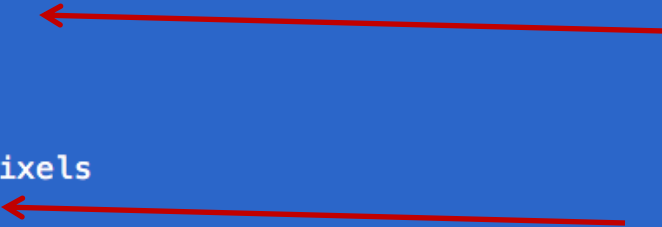


Optimal data collection ?



Header – is CBF header a MAH ?

```
# 2015/May/06 10:30:40
# Pixel_size 172e-6 m x 172e-6 m
# Silicon sensor, thickness 0.001 m
# Oscillation_axis omega
# Excluded_pixels: badpix_mask.tif
# Chi 0.0000 deg.
# Angle_increment 0.1000 deg.
# Polarization 0.99
# file_comments
# N_oscillations 2500
# Beam_xy (1223.03, 1256.56) pixels
# Exposure_time 0.020000 s
# Phi 0.0020 deg.
# Energy_range (0, 0) eV
# Start_angle 160.6000 deg.
# Detector_distance 0.617619 m
# Detector_Voffset 0.0000 m
# Alpha 0.0000 deg.
# Flat_field: (nil)
# Threshold_setting 7619 eV
# Exposure_period 0.020950 s
# N_excluded_pixels: = 321
# Kappa 0.0020 deg.
# Tau = 0 s
```

Two red arrows originate from the right side of the image. The first arrow points to the line '# Angle_increment 0.1000 deg.' and the second arrow points to the line '# Exposure_time 0.020000 s'.

Do you like this image?



Do you like this image ?



How expensive is bright lens ?



[See more choices](#)

Canon EF 85mm f1.2L II USM Lens
for Canon DSLR Cameras - Fixed
by Canon

\$1,999.00 ✓Prime

Get it by **Monday, Aug 24**

More Buying Choices

\$1,999.00 new (22 offers)

\$1,499.99 used (24 offers)

Trade-in eligible for an Amazon gift card

★★★★★ ▾ 159



[See Style Options](#)

Canon EF 85mm f/1.8 USM Medium
Telephoto Lens for Canon SLR
Cameras - Fixed
by Canon

\$369.00 ✓Prime

Get it by **Monday, Aug 24**

More Buying Choices

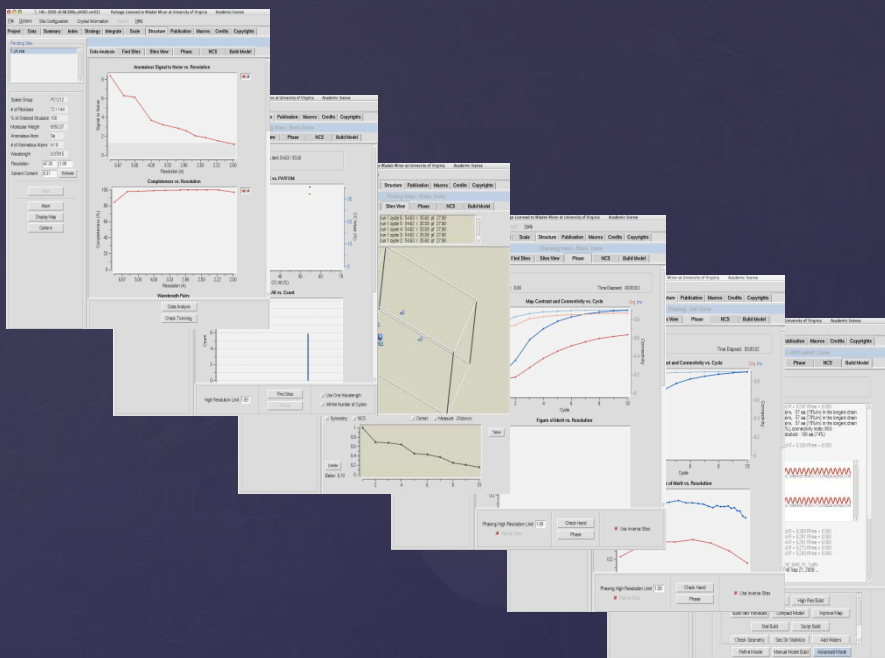
\$369.00 new (27 offers)

\$298.00 used (26 offers)

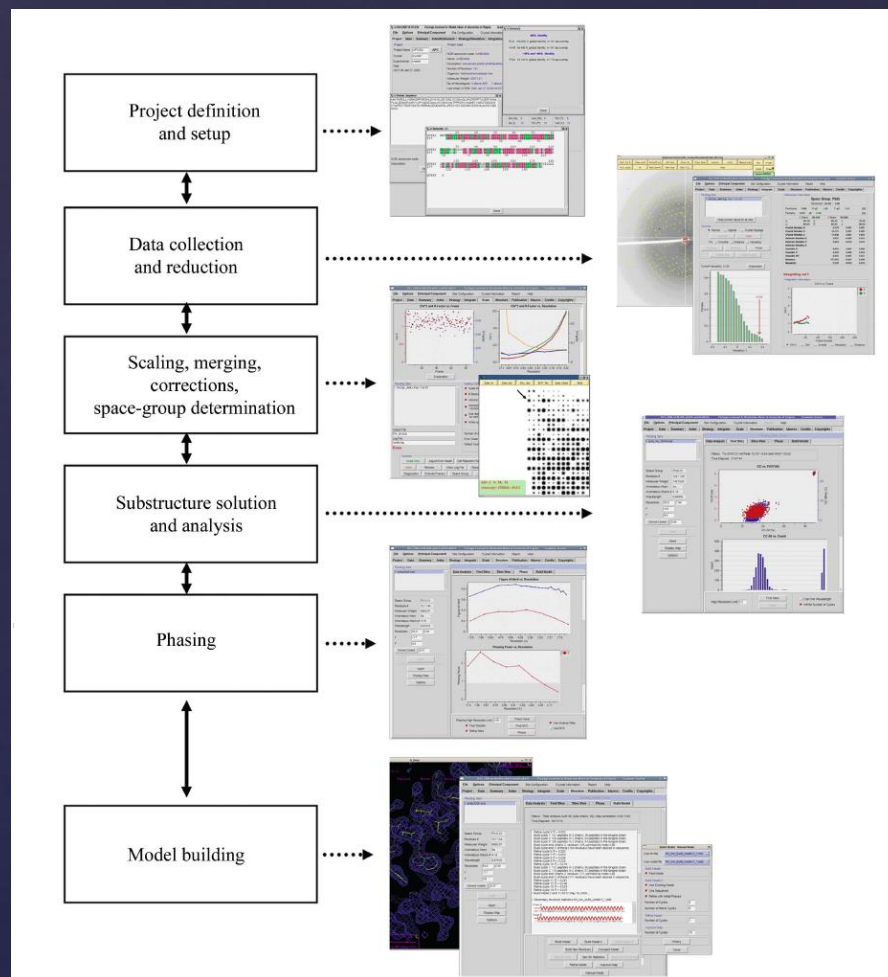
Trade-in eligible for an Amazon gift card

★★★★★ ▾ 770

HKL-3000 (6 mouse clicks program)



SHELXD, SHELXE
CCP4, DM, REFMAC
SOLVE, RESOLVE
ARP/WARP
O, COOT, CCP4

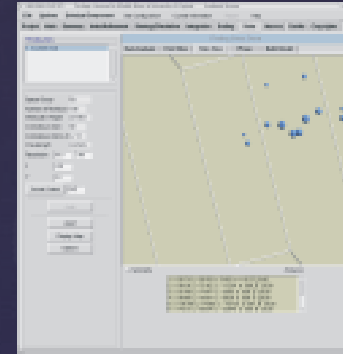
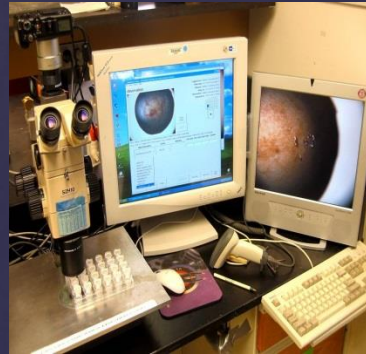


HKL-3000 at SBC

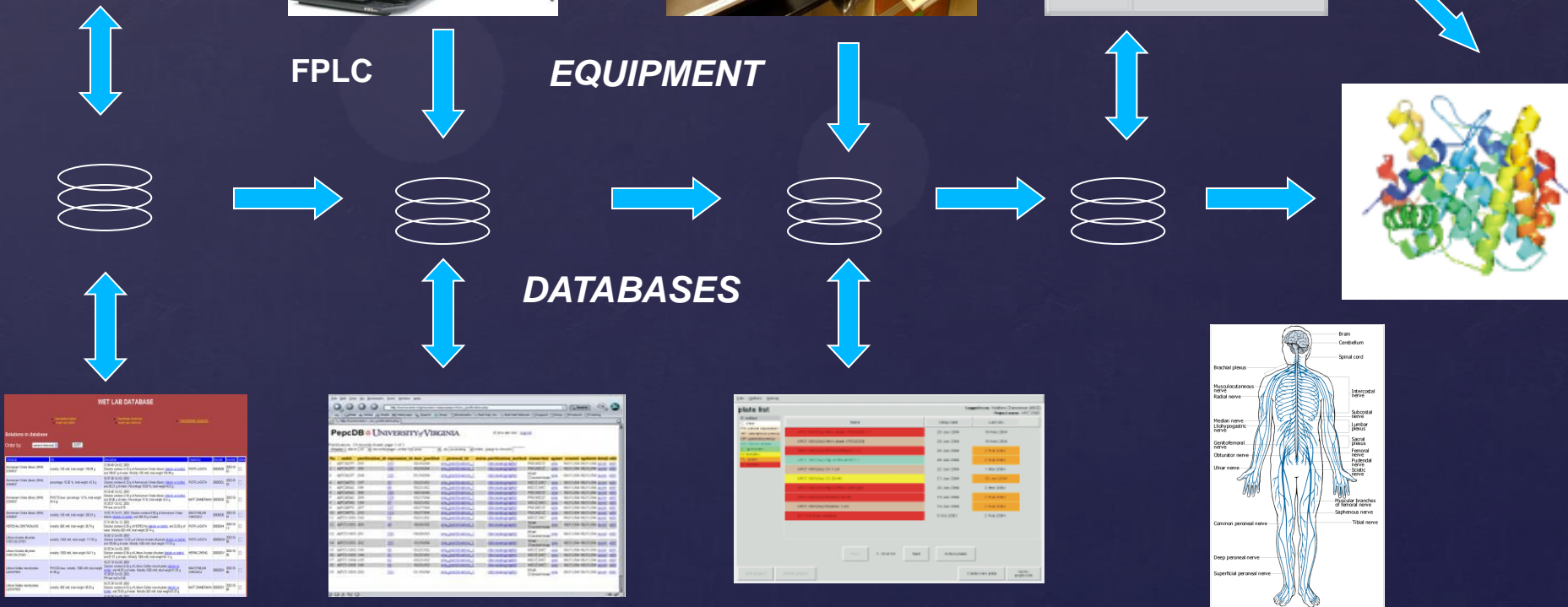


Database-controlled pipeline

lab e-book



HKL-3000



Big brother?

Statistics / Progress in Minor Lab LIMS by researcher

Last week (17 Apr 2015 - 24 Apr 2015)

Person	Clones	Exprs	Purifs	Macro preps	Plates	Drops	Crystals	Datasets processed	Structure refs	Kinetic assays	Thermal shift assays
Cooper, David	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>23</u>	<u>18</u>	<u>0</u>	0	0
Handing, Katarzyna	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>51</u>	<u>53</u>	<u>13</u>	0	0
Hou, Jing	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>1</u>	30	<u>0</u>	<u>1</u>	<u>1</u>	0	0
Kowiel, Marcin	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>1</u>	<u>8</u>	<u>3</u>	0	0
Shabalin, Ivan	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>125</u>	<u>14</u>	<u>9</u>	0	0
Shumilin, Igor	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>0</u>	<u>3</u>	<u>2</u>	0	0
Szlachta, Karol	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	0	<u>34</u>	<u>20</u>	<u>3</u>	0	0

Last month (25 Mar 2015 - 24 Apr 2015)

Acknowledgments

Wladek Minor

- Matt Zimmerman
- Marek Grabowski
- Heping Zheng
- Marcin Cymborowski
- Karol Langner (Google)
- Przemek Porebski
- Piotr Sroka
- Ivan Shabalin
- Katherine Handing

Zbyszek Otwinowski

- Dominika Borek

Andrzej Joachimiak

MCSG and SBC staff

Wayne Anderson and CSGID staff

Steve Almo and NYSGRC Staff

Ian Wilson, Marc Elsliger and JCSG staff

Steven Burley, John Westbrook and PDB staff

Tom Terwilliger

Zbyszek Dauter

Grants:

U01-HG008424

NIH GM53163, GM62414, GM74942

GM093342, GM094585, GM094662

DOE, NCI

NIAID HHSN272200700058C

NIAID HHSN272201200026C

HKL Research. Inc.