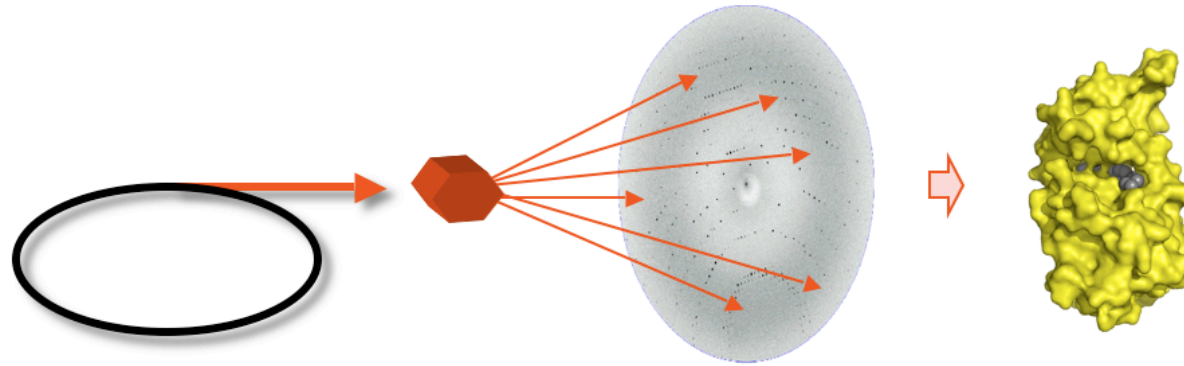# Data in Protein Crystallography

Thomas R. Schneider
Project Coordinator EMBL@PETRA3

Standard Data Formats for Experiments with
Photons, Neutrons, Ions
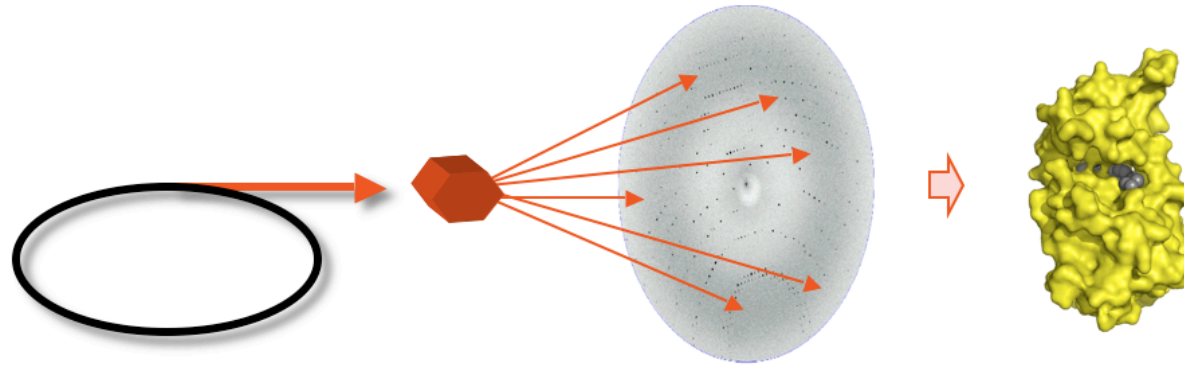DESY 27/10/2010

EMBL

# Macromolecular Crystallography (MX)



- **Raw Data:** 2D diffraction images collected from a sample rotating slowly around a chosen axis.

- **Extra:** Many crystals have to be measured(100-1000ds) potentially spatially resolved; compatible data sets will be combined later on.

- **Boundary conditions:**
  - Users use various synchrotrons
  - Users are not physicists
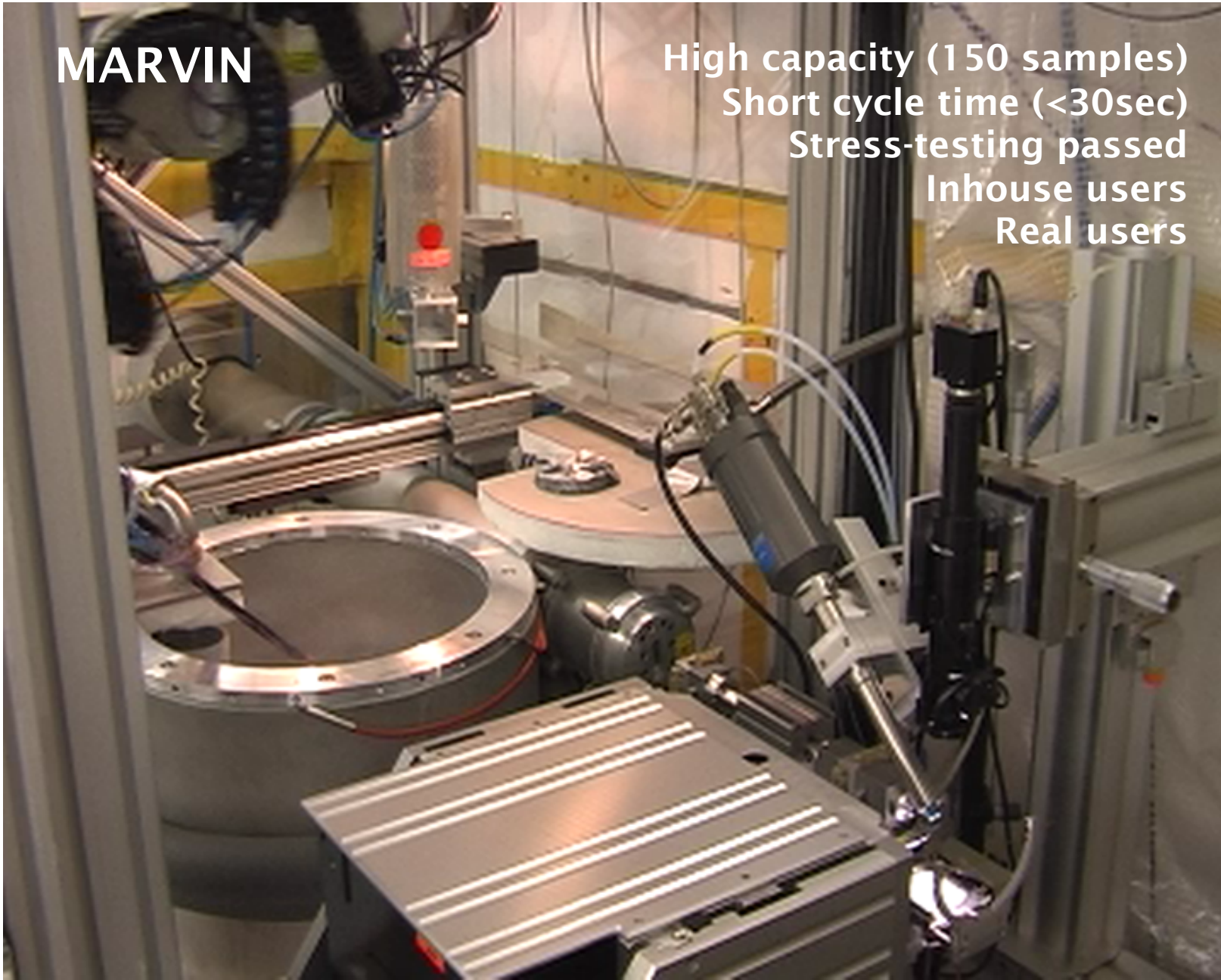  - 'Upstream' and 'Downstream' data are very important as well

EMBL

# The Plan



1. **Procedures producing data**
2. **Some approaches in use in MX**
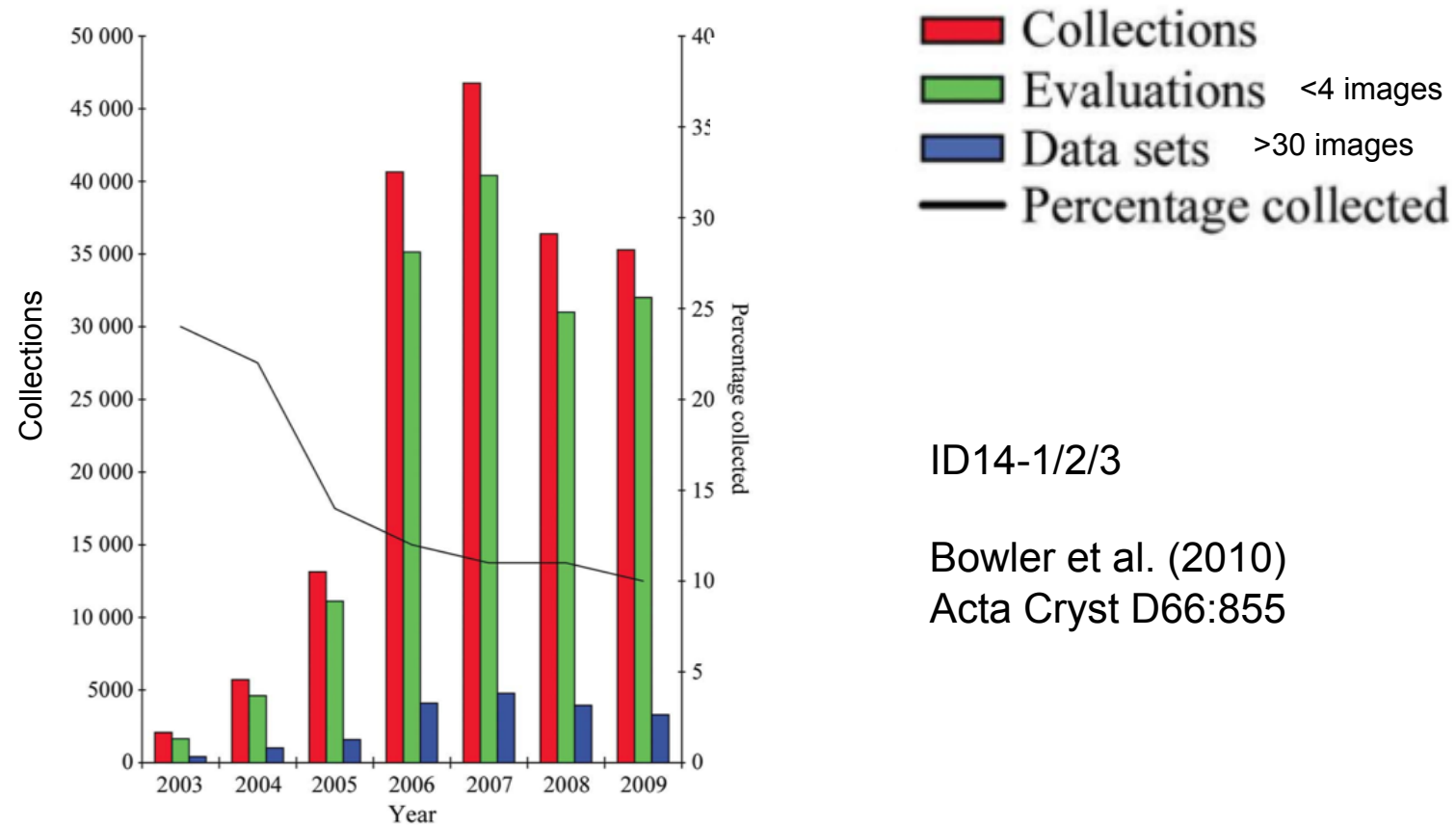
EMBL

# Raw image data

- A typical data set collected on a PILATUS 6M detector using the shutter-less rotation method can be collected in **2.5 min.** and produces ca. **4.5 GB** of raw data. 'In phase' processing is possible using local clusters.

- On a well equipped beamline, in principle, 250 data sets can be collected in 24 hours. -> **1 TB of raw data / day.**

- New technologies are already / will be soon available.
  - Frame-shift CCDs. Frame rate -> 100 Hz
  - Next generation PADs. Frame rate -> kHz, more pixels
  - Next generation amorphous Se detectors

- On 3$^{rd}$ generation sources, MX is not flux-limited -> more samples per time are possible.

- In principle there is CBF (Crystallographic Binary File Format, http://www.esrf.eu/computing/Forum/imgCIF/cbf_definition.html, cif-type header + binary image), however it can be difficult to motivate detector manufacturers to obey the rules (abuse of free format fields …).

EMBL

**MARVIN**

High capacity (150 samples)
Short cycle time (<30sec)
Stress-testing passed
Inhouse users
Real users

EMBL

# Evaluation vs. Data Collection



ID14-1/2/3

Bowler et al. (2010)
Acta Cryst D66:855

EMBL

# Sheer Volume

- http://www.esrf.eu / UsersAndScience/ Experiments/MX
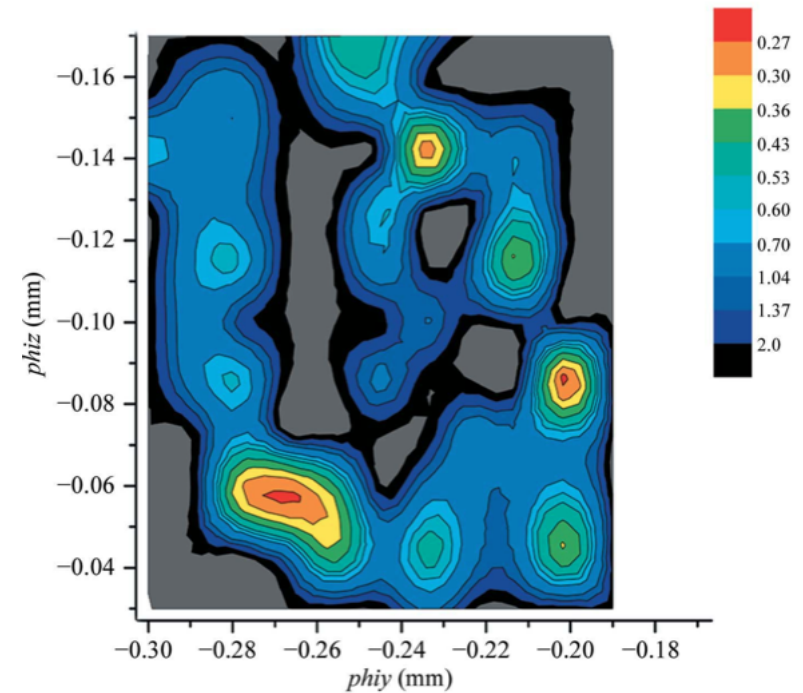
**Data collections 2010**

Wed 27 Oct, 09:17

- **id14eh1** Sample Evaluations: 9921, Data Sets: 1376

- **id14eh2** Sample Evaluations: 3799, Data Sets: 356

- **id14eh4** Sample Evaluations: 15137, Data Sets: 3196

- **id29** Sample Evaluations: 15639, Data Sets: 4320

- **id23eh1** Sample Evaluations: 22371, Data Sets: 4528

- **id23eh2** Sample Evaluations: 21885, Data Sets: 2572

- **Total Sample Evaluations: 88752, Total Data Sets: 16348**

ISPyB

EMBL

# Fine Grained Sample Evaluation



Beta1-andregenic GPCR

Rastering in 15 micron steps + scoring

- 'Diffraction Cartograph' Bowler et al. Acta Cryst D (2010) D66:855

EMBL

# Data Collection Strategy

- Results from sample characterization have to be analysed and a strategy has to be created to obtain the best possible data for the purpose (phasing, ligand structures ….) while fighting radiation damage

- BEST (Gleb Bourenkov, EMBL-HH; Sasha Popov, ESRF-GR)

- Conditions for 'clever' data collection

    - Evaluation parameters and resulting data are stored

    - Crystals can be remounted reproducibly (within X micron, within Y degree) -> 'NewPin'

- The 'real' data collection could actually take place on a different beamline / synchrotron

EMBL

# Workflows / Data Models / EDNA

- Framework for implementing pipelines for online-data analysis.

- At the heart is the data model for the specific process. This can be defined e.g. in XML schema definition (XSD).

- Establishing a data model for kappa-goniometry with one flat detector has been a difficult exercise.

- EMBL-HH has written an mxcube-brick for integrating EDNA into the beamline user interface.



Incardorna et al. (2009)
J. Sync. Rad. 16:872

EMBL

# Sample Management is needed for 100/1000s of samples

EMBL

# The SPINE standard sample holder



Cipriani et al. (2006)
Acta Cryst. D62:1251

EMBL

# Sample Management



- Items
    - Samples (barcode, protein acronym, crystal form)
    - Containers / baskets (barcode)
    - Dewars (barcode, courier tracking number)
    - Shipments (set of dewars)
    - Pile of shipments (Fri @ ESRF)
- Physical location
    - Where is the crystal?
    - Where is the dewar…

http://www.esrf.eu/UsersAndScience/
Experiments/MX/
How_to_use_our_beamlines/ISPYB/
ISPyB_090915_01%20_00.pdf

EMBL

# iSpyB



ESRF Experiment Division Tuesday Events
15/09/2009

**ISPyB**
Information System for Protein crYstallography Beamlines

From your sample to your data analysis: how to track every step of your experiment in a database. An example with ISPyB for MX experiments

Patrice Brenchereau ESRF/CS/MIS

European Synchrotron Radiation Facility

# iSpyB – Data Collection

# iSpyB – Data Collection

# iSpyB – Data Mining

EMBL

# iSpyB – Sample Ranking

# iSpyB – Data Model



SHIPMENT / DEWAR / CONTAINER

SAMPLE / PROTEIN / CRYSTAL

EXPERIMENT

DATA COLLECTION

EDNA

EMBL

# iSpyB is open source

# Data Evaluation

- 'Standard cases'
  - HKL2MAP will produce an electron density map in 3-5 min.
  - AutoRickshaw will produce a 3D-model in some hours

- Difficult cases
  - Many dead ends

- Inclusion of HKL2MAP results into iSpyB has been attempted.

EMBL

# TARDIS

- http://tardis.edu.au
- http://
  www.monash.edu.a
  u/news/newsline/
  story/1608, http://
  mpegmedia.abc.net.
  au/rn/podcast/
  2010/07/
  fte_20100715_0850
  .mp3

# Structural Genomics projects

- Go from gene to structure
- Well-defined workflows including experiments at synchrotrons
- Examples:
  - Joint Center for Structural Genomics www.jcsg.org
  - Structural Genomics Consortium www.thesgc.org



http://www.nigms.nih.gov/Initiatives/PSI/Centers/JCSG.htm

EMBL

# Long-Term Archiving of raw data

- Reasons
  - Potential value for re-evaluation with new technology / fresh brain
  - 10 years obligatory documentation
  - Traceability in case of possible fraud

- Who owns the raw data / meta data?
  - The PI? For how long?
  - The synchrotron?
  - The public?

- Role management
  - Assignment of roles to users
  - Identification of users and their roles
  - mobility between groups/institutions
  - orphaned data (retirement etc.)

- dCache@DESY plus an interface (Frank Schluenzen, Ilya Agapov)

EMBL

# Some additional points

- 'Religious' approach to work flows (different labs have different religions) -> flexibility / intuitivity

- Confidentiality / Visibility / Practicality issues with centralized servers.

- Practical work with users:
  - Much of the practical work is done by inexperienced PhD students
  - More experienced scientists only sporadically collect data (at different sources)
  - If meta-data collection is not done automatically, it will not be done (in academia at least)
  - Data Management should be totally transparent (drag-and-drop …)

- Licensing issues with data handling systems (GPL, LGPL, open source, unknown, …)

EMBL

# EMBL@PETRA3

- **Goal**: 3D structural information about biological macromolecules and their complexes at the highest possible resolution (see ribosome).

- **'Integrated Facility for Structural Biology'** will offer user access to:

  - 1 SAXS (small angle x-ray scattering) beamline at PETRA III

  - 2 MX (macromolecular crystallography) beamlines at PETRA III

  - High-Throughput Crystallization Facility (1 Mio experiments, largest in Europe)

  - Sample preparation and characterization

  - (Remote) Data Evaluation

    - arpWarp, ATSAS, autoRickshaw, …

- Remote Access where requested and possible

- Industrie is welcome – confidentiality is important

EMBL

# MX – Different data at different steps

- Protein Production -> 'PIMS' et al. (>2nd inc, >10y)

- Protein Characterization -> 'PIMS' et al.

- SAXS experiment -> 'iSpyB' (nascent)

- Crystallization -> 'Crims' et al. (>5y)

- Crystallographic Testing and Data Collection -> 'ispyB' (2nd inc, >5y)

- Data Processing -> 'ispyB'

- Data Evaluation -> 'ispyB'

- Results -> Protein Data Bank (+20y)

- Long time Archive -> ???

EMBL

# Summary

- Synchrotron-based crystallography is an advanced technique with established work-flows.

- Despite this, in most places in Europe (except ESRF and DIAMOND) data management is still done by transporting hard disks.

- This system will not be viable in the future:

  - Robotics will further increase through-put

  - Fast detectors will increase through-put and data volume

  - Requirements for documentation will be enforced

- Data from related experiments / methods (also wet-lab) need to be stored and managed ideally in an integrated manner. Consistent meta data / connectivity between data are very important.

- Resources needed for design and implementation are substantial

- **Working together with other facilities is important both from the user side (homogeneity) and the supplier side (synergies).**

EMBL