

EDNA

A tool for online data analysis

and

A collaborative effort

Olof Svensson

EDNA Project Manager

ESRF Scientific Software Group

Overview

- Online data analysis
 - Example: MX automation
- The EDNA prototype
 - Data collection taking into account radiation damage
 - Results
- The EDNA Project / Framework
 - Project management
 - Data model
 - Modularity / plugins
 - The testing framework
- Future developments
 - The first EDNA release (Kernel + MXV1): May 2008
 - MXV2 / Tomography

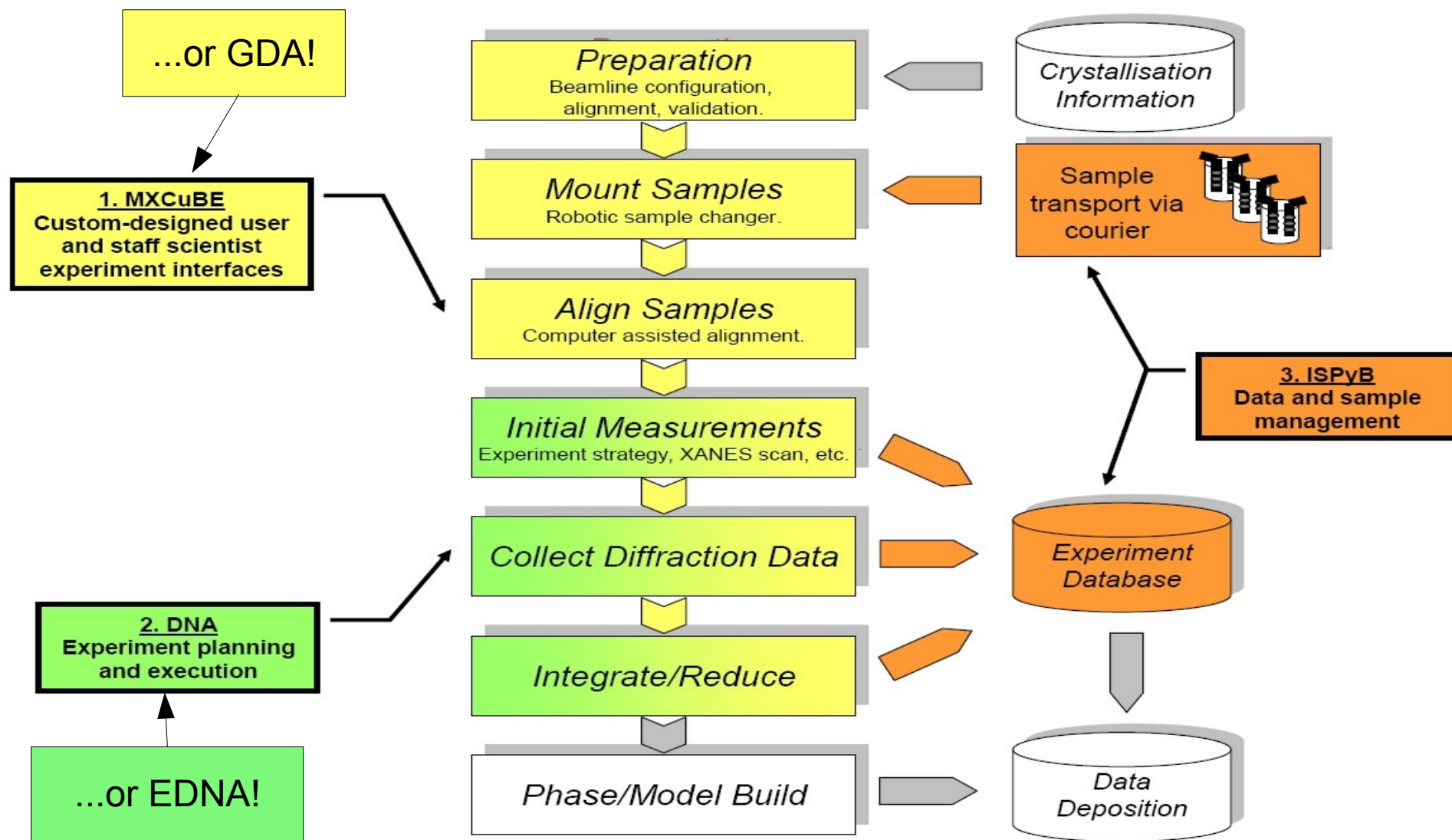
Online Data Analysis (1)

- Automation
 - Example MX : Data Collection Pipeline (DCP)
- Grid / Batch processing
 - Fast data processing \Rightarrow parallel / distributed execution

Online Data Analysis (2)

- Data Management
 - Lims
- Modular
 - It must be easy to change the work flow / scientific programs
 - Data Model
- Robust - code reliability
 - If ODA fails once, the user won't try again...
 - Efficient error tracking system
- Collaboration
 - Avoid re-inventions of the wheel
 - Project management
 - Documentation, Use Cases

MX Automation - the Data Collection Pipeline: Full automation from sample loading to reduced (integrated and scaled) data

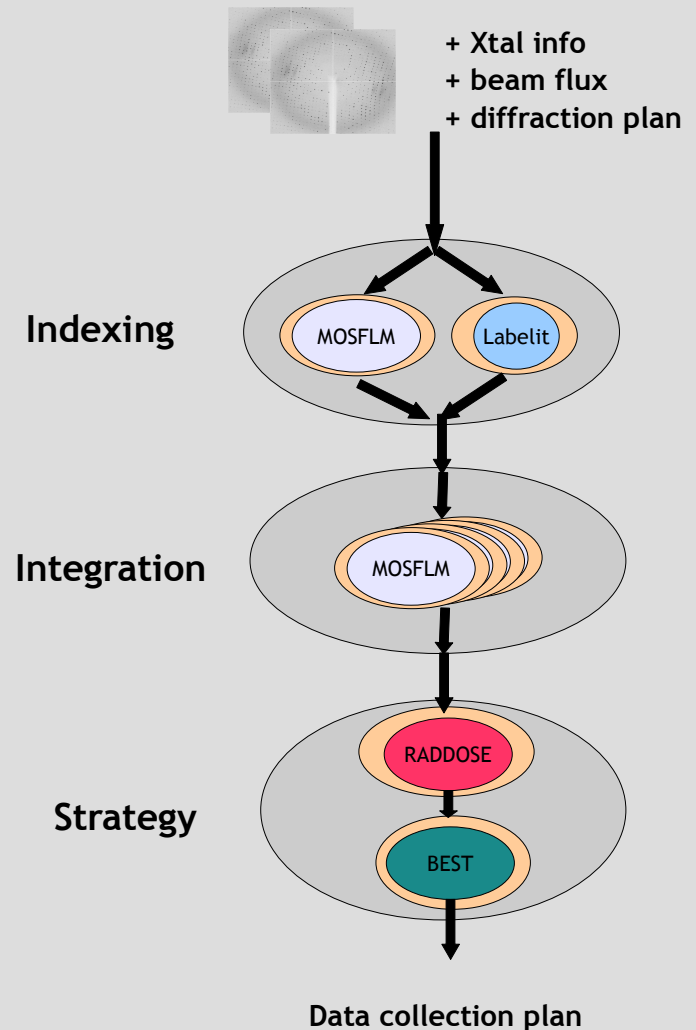


What is EDNA?

- EDNA is an international collaborative project between several institutes and synchrotron facilities.
- Developed on the foundation of the project automated collection of data (« DNA », www.dna.ac.uk)
- Designed to be a framework for Online Data Analysis of X-ray experiments

The EDNA Prototype

- MX sample characterisation taking into account radiation damage
 - Indexing using MOSFLM or Labelit
 - Parallel integration of reference images
 - If flux + beamsize + chemical composition: RADDOSE for estimating radiation damage
 - BEST strategy calculation
 - taking into account radiation damage
 - multi-subwedge data collection strategies



NSLS CBASS GUI

The screenshot displays the NSLS CBASS GUI interface. At the top, there's a menu bar (File, Edit, Tools, Comments) and a toolbar with buttons for Collect, Setup, Pucks, and Canes. Below this is a table for experiment parameters:

	Start	End	Width	Time	File Prefix	NumStart	Distance	Wavelength
1	65.0	113.3	1.05	43.002	x296	0	543.05	1.1000
2								
3								
4								

Below the table are control panels for Goniostat (Omega, Kappa, Phi) and Detector Status (Distance, Edge Resolution, Tilt). A circular diagram on the right shows detector positions at 1.90 Å and 2.52 Å. A row of buttons includes Collect Data, Collect Strategy (EDNA), Open Shutter, Close Shutter, Home, Pause On Beam Dump, Clear Collection, Clear Top Row, Pause, Count, and Abort.

The status bar shows the file path: /home/pxuser/cbass_test/ref-x296_1_002.img and a green "CBASS Ready" indicator.

A table of experimental results is shown below the status bar:

Command:	Group...	skinner	PxID...	px04-4010
Strategy : Best : 7.90 7.25 97.0 1976.3 119.6 16.5 15.0 1.21 5.0 0.00				
Strategy : Best : 7.25 6.73 96.0 1264.8 121.9 10.4 9.8 1.11 7.6 0.00				
Strategy : Best : 6.73 6.31 99.0 928.0 130.7 7.1 6.9 1.07 10.6 0.00				
Strategy : Best : 6.31 5.96 100.0 736.6 142.2 5.2 5.1 1.05 14.0 0.00				
Strategy : Best : 5.96 5.67 99.0 616.1 148.8 4.1 4.1 1.04 17.6 0.00				
Strategy : Best : 5.67 5.41 99.0 546.1 161.0 3.4 3.3 1.03 21.2 0.00				
Strategy : Best : 5.41 5.19 100.0 515.3 169.1 3.0 3.0 1.03 23.6 0.00				
Strategy : Best : 5.19 4.99 100.0 481.2 182.5 2.6 2.6 1.02 27.1 0.00				
Strategy : Best : 4.99 4.81 100.0 441.2 188.7 2.3 2.3 1.02 30.8 0.00				
Strategy : Best : 4.81 4.65 100.0 399.8 199.9 2.0 2.0 1.02 36.3 0.00				
Strategy : Best : 12.00 4.65 98.0 1265.6 160.0 7.9 7.6 1.09 8.8 0.00				

An "EDNA Input Parameter" dialog box is open, showing values for Completeness (0.9900), I/Sigma (2.0000), Multiplicity (auto), and Resolution (auto). Buttons for Apply and Close are visible.

Below the table, a text area shows the following information:

```

Strategy
EDApplicationPrototype-v1.0.0::Runtime: 29.2395579815 [s]
spacegroup = P3 mosaicity = 0.75 resolutionHigh = 4.65 cell_a = 86.009 cell_b = 86.009 cell_c = 156.4334 cell_alpha = 90.0 cell_beta = 90.0 cell_gamma = 120.0 status = ok
dna Strategy results: Start=65.0 End=113.3 Width=1.05 Time=43.002 Dist=543.05
collection_dist 543.05
CBASS>
    
```

At the bottom, a table summarizes the experiment details:

group	px id	xtal id	file template	sweep	# images	tot exp. time	timestamp
skinner	px04-4010	Standard Project					
	x296	x296		0	2	10	27-OCT-2008 10:45

Below this table, a detailed summary of the experiment parameters is provided:

spacegroup = P3 mosaicity = 0.75 resolutionHigh = 4.65 cell_a = 86.009 cell_b = 86.009 cell_c = 156.4334 cell_alpha = 90.0 cell_beta = 90.0 cell_gamma = 120.0 status = ok dna Strategy results: Start=65.0 End=113.3 Width=1.05 Time=43.002 Dist=543.05

The bottom of the screen shows a taskbar with various application windows open, including [Termin..., [EDNA ..., [firewir..., [pxuse..., [xterm], [pxuser..., [pxsys..., [xterm], [pxuser..., [pxuse..., and [xterm].

EDNA Prototype ccp4i GUI

Title EDNA prototype

Run EDNA Characterization using **Images** input

☒ Account for radiation damage using chemical composition input

Data set # 1

image #1 in **FAE** Browse View

Edit list Add an image

glob clean

Edit list Add data set

XML Output **FAE** Browse View

Diffraction Plan

☐ Force Space Group

Strategy complexity: few sub-wdges

Maximum Exposure time per data collection 100000 seconds

Aimed I Over Sigma at highest resolution 3.0

Minimum oscillation width 0.2 degree(s)

☐ Define Aimed Completeness (default >= 0.99)

☐ Define Aimed Resolution (default - highest possible)

☐ Define Aimed Multiplicity (default - optimized)

Beam

X-ray beam

Flux x 10¹¹ photons/second

Beam size (Horizontal x Verical) x mm²

Sample

Dimensions 0.1 x 0.1 x 0.1 mm³ Radiation susceptibiliy 1

Chemical composition

Solvent

Run Save or Restore Close

Program author: EDNA developers <http://www.edna-site.org>

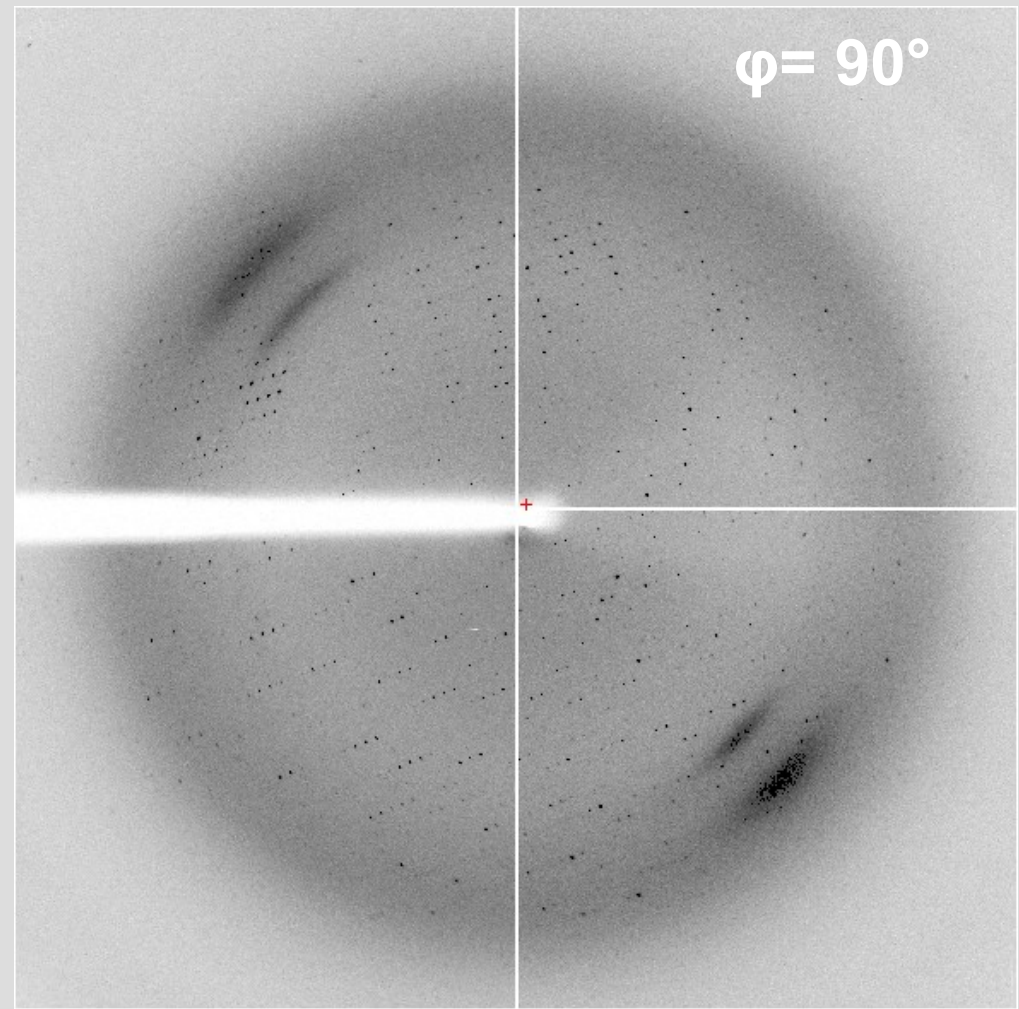
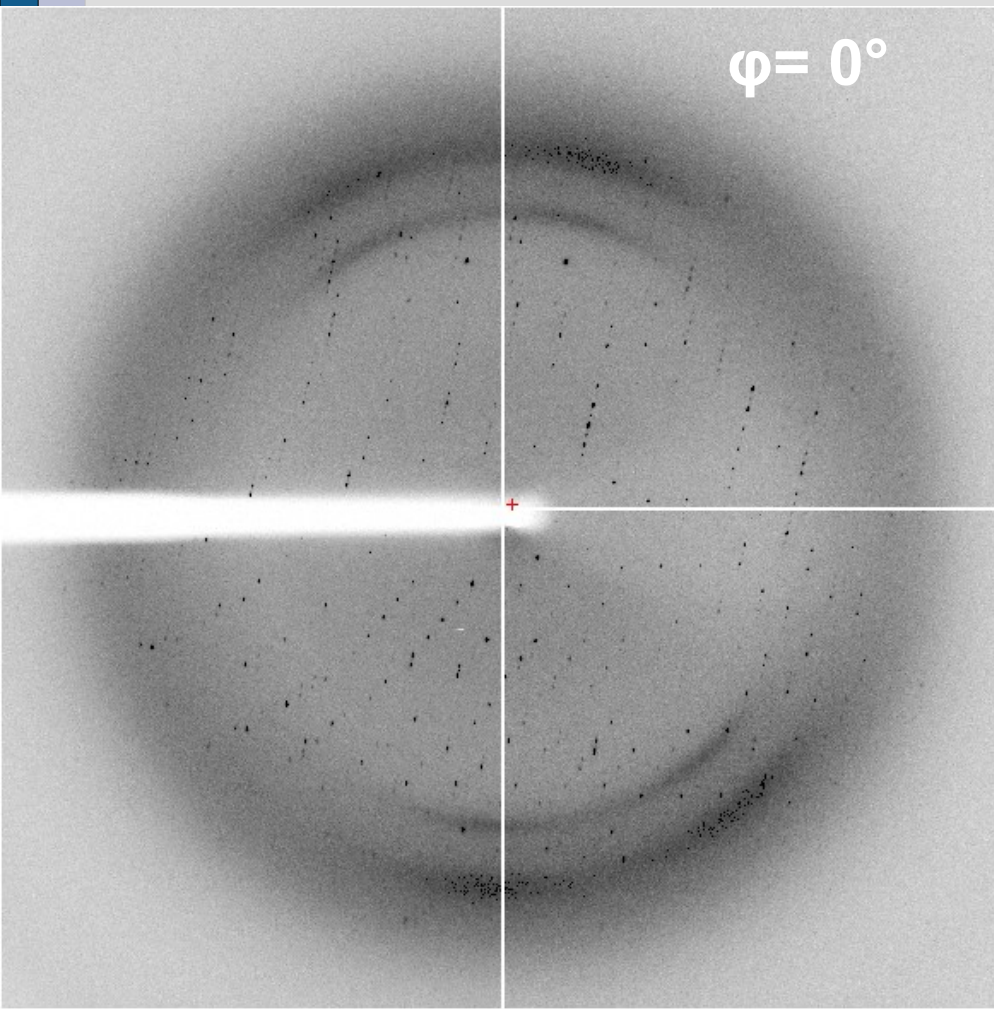
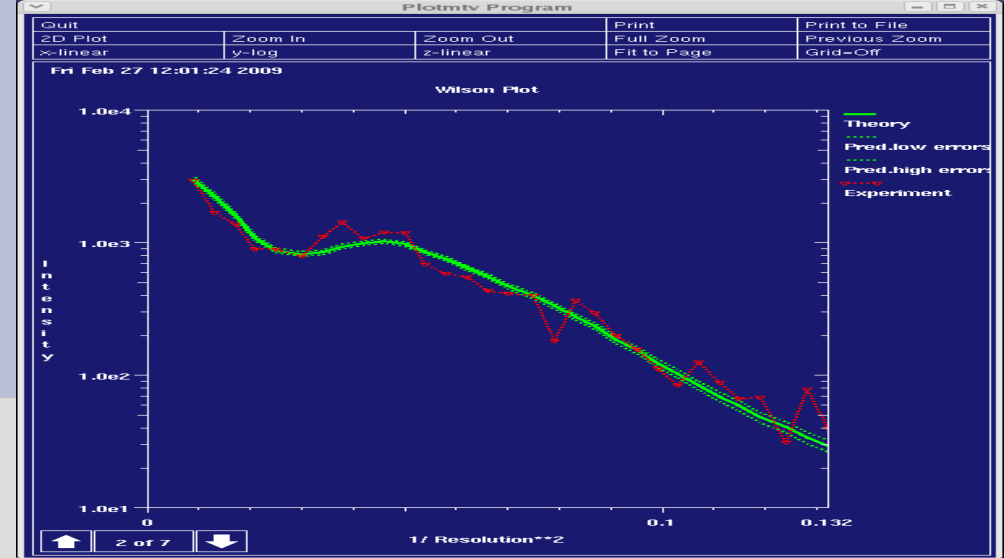
ID14-2, ESRF

Initial images for Survival E protein crystal

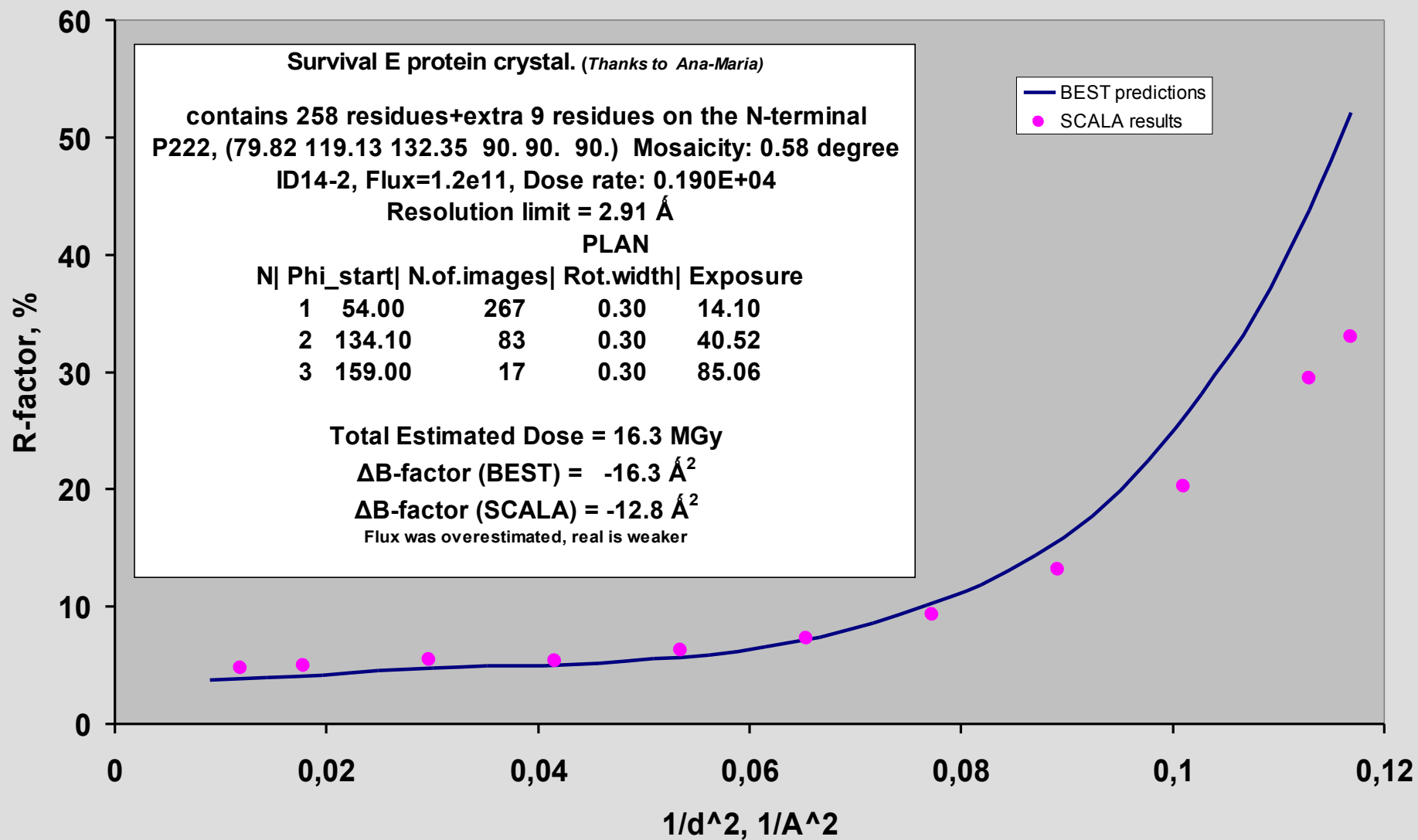
Exposure = 5 s,

Rot. width = 1°,

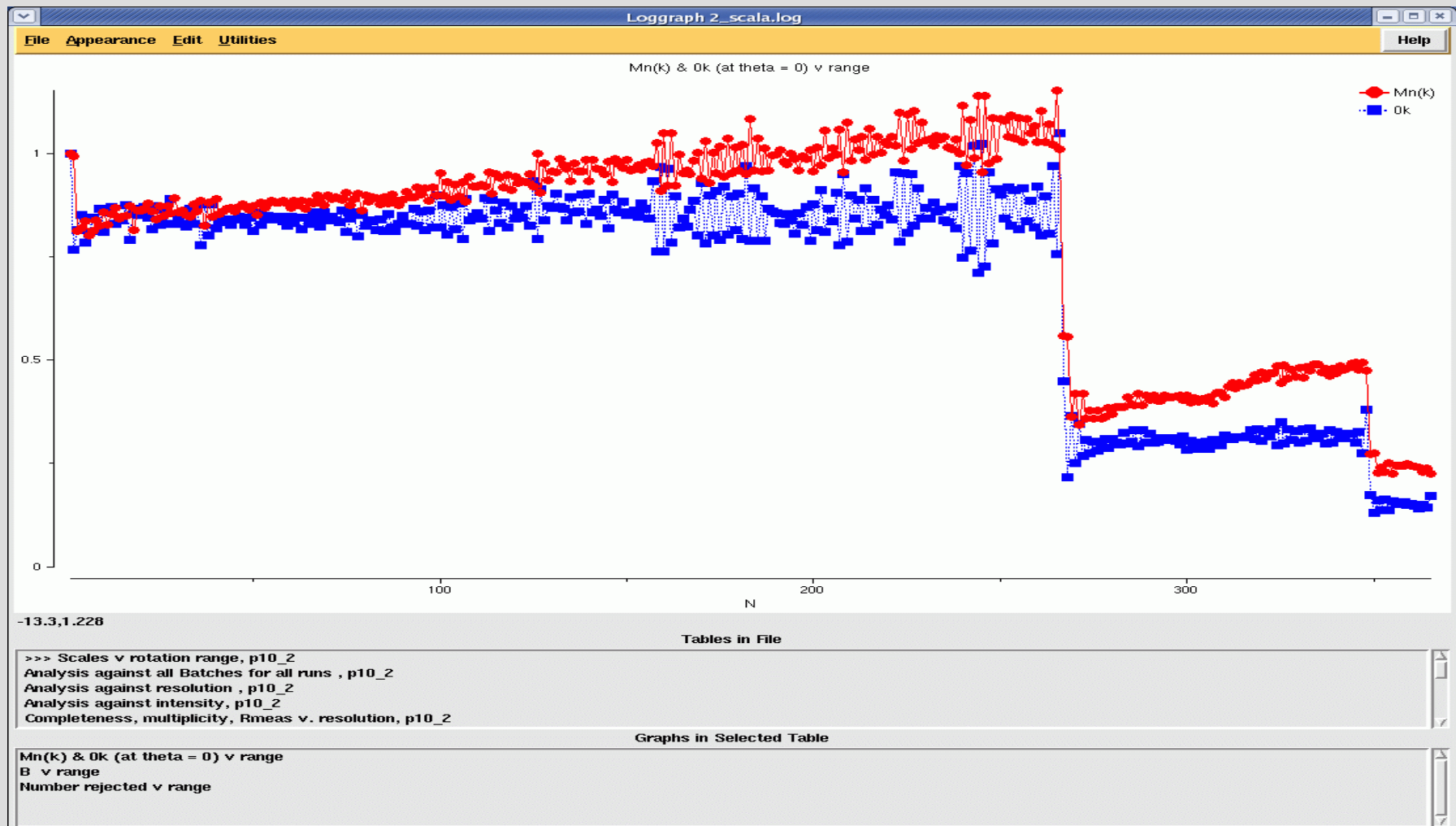
Resolution=2.8 Å



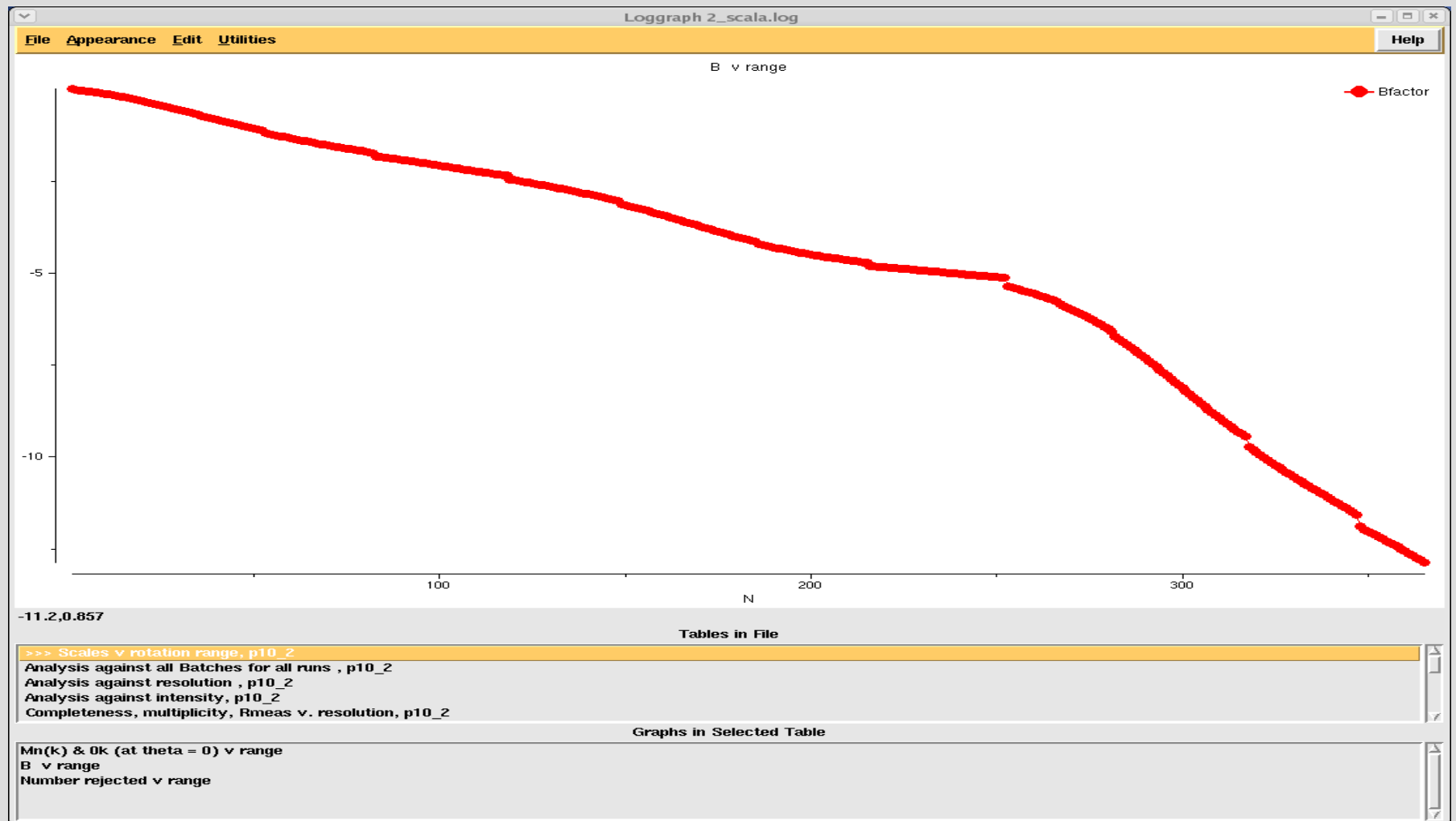
R-factor VS. Resolution



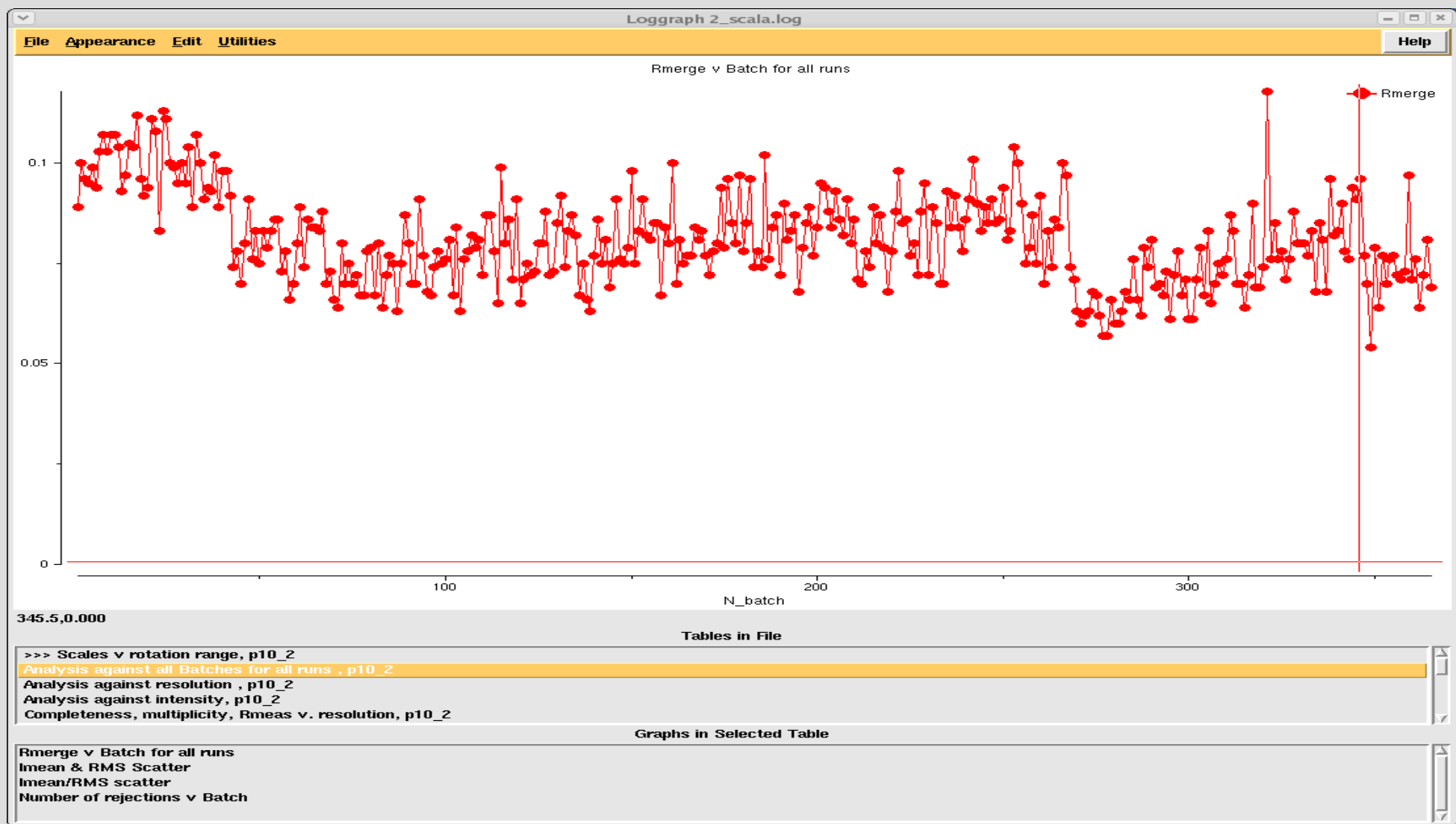
SCALA Output



SCALA Output

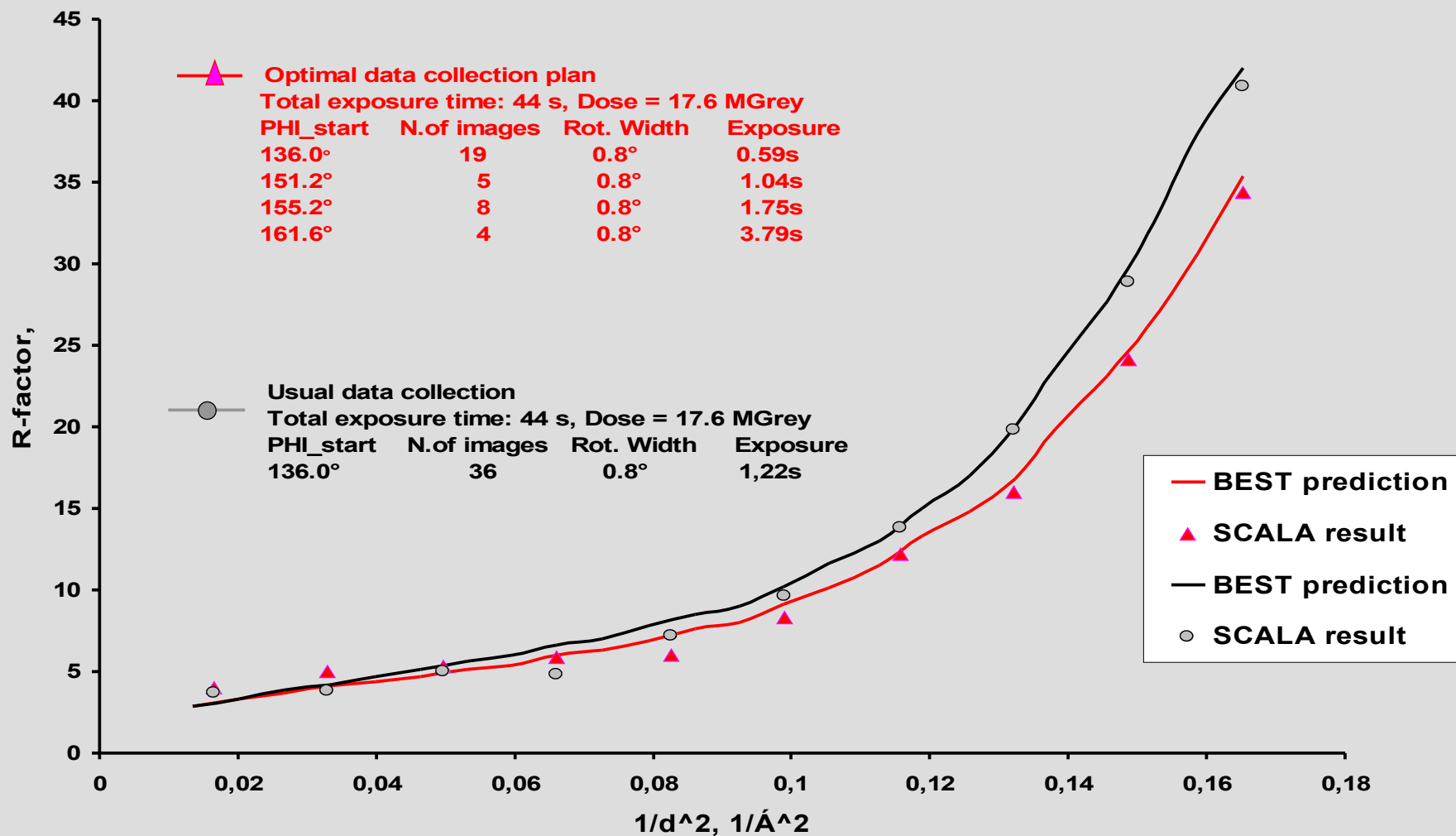


SCALA Output



ID23-1, ESRF

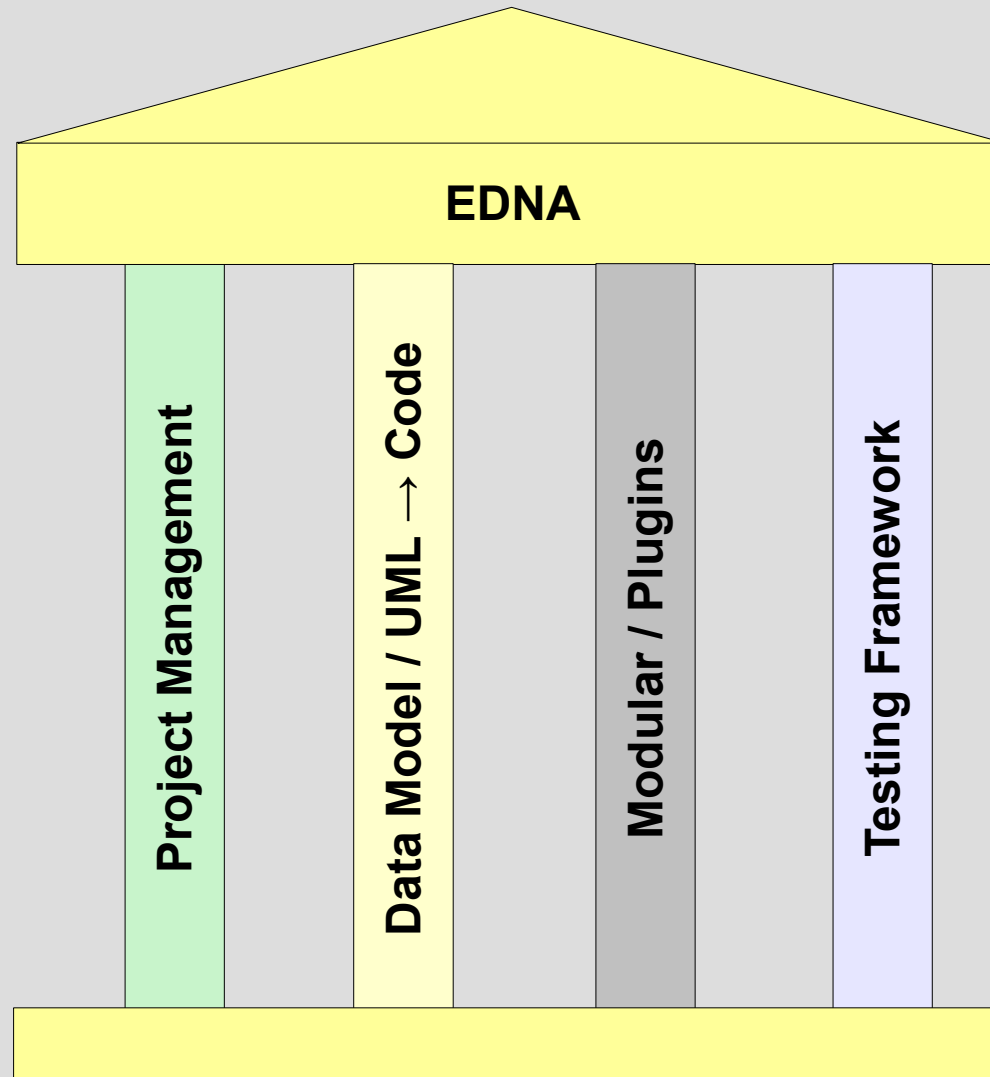
Complex of viral (from tomato bushy stunt virus, TBSV) supressor of RNA silencing p19 with mismatched siRNA



EDNA is not DNA...

- Same goals (initially)...
 - Automatic MX sample characterisation
 - Online data processing during MX data collection
 - Ranking
- ...however very different implementations
 - EDNA designed to not be specific to MX
 - No shared code base between DNA and EDNA
 - Different project management
 - Different collaborators
- If time allows:
 - Short presentation of DNA
 - History behind DNA → EDNA

The EDNA Project / Framework



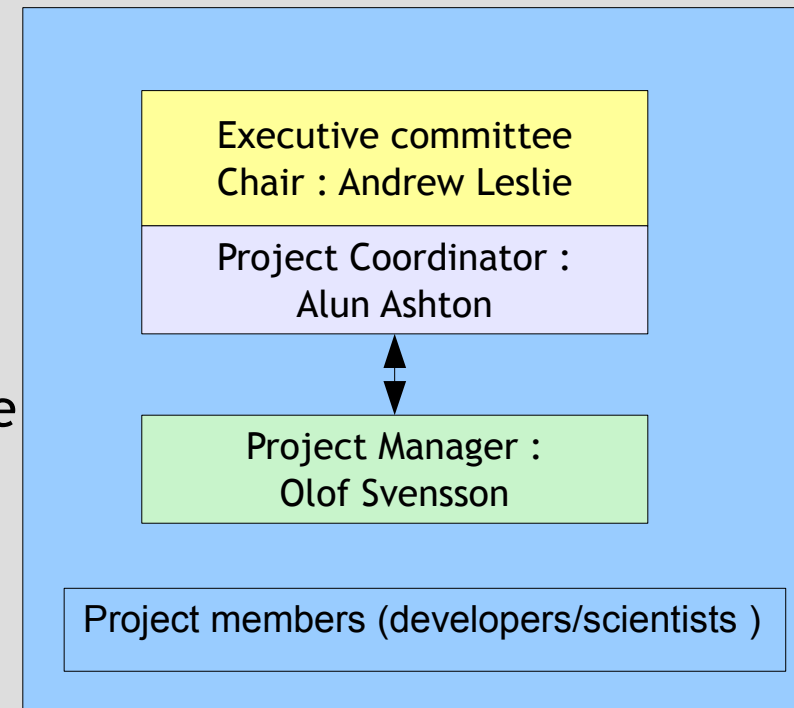
EDNA Project Management (1)

- Executive Committee:

- Alun Ashton, DLS, UK
- Gérard Bricogne, Global Phasing, UK
- Andrew Leslie, MRC LMB, Cambridge, UK
- Andrew McCarthy, EMBL-Grenoble, France
- Sean McSweeney, ESRF, Grenoble, France
- Thomas Schneider, EMBL-Hamburg, Germany
- Andrew Thompson, Synchrotron Soleil, France

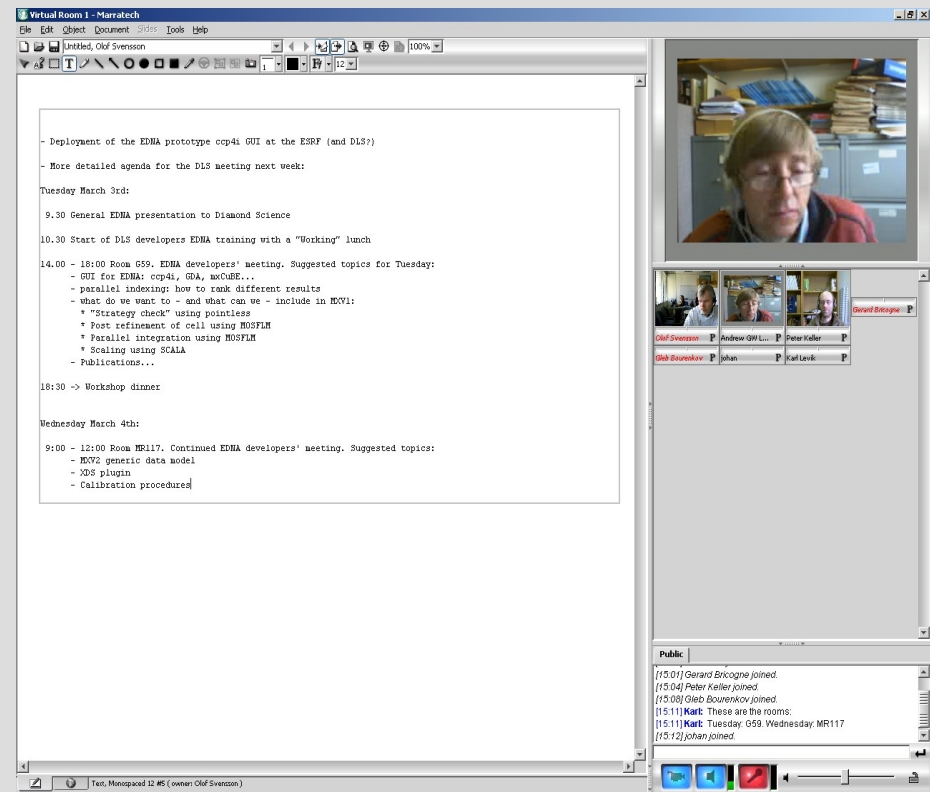
- Other members from:

- BESSY, Berlin, Germany
- MAX LAB, Lund, Sweden
- NSLS, Brookhaven, U.S.
- SLS, Villigen, Switzerland
- University of Sydney, Australia
- University of York, UK



EDNA Project Management (2)

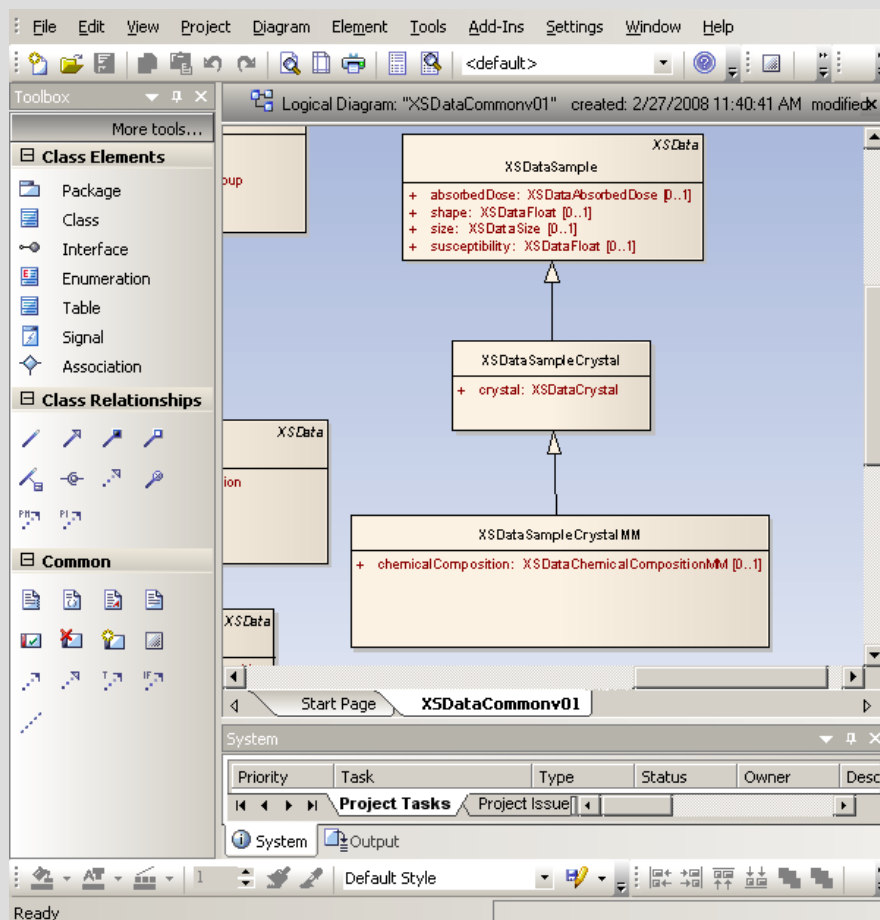
- Project agreement
- Coding conventions
- Code reviews
- Development tools
 - Eclipse
 - Enterprise architect
- Project portal
 - <http://www.edna-site.org>
 - Wiki documention
 - Bugzilla server
 - Subversion server
 - Discussion forum
- Executive committee
- Video conferences
- Developers' meetings & workshops



Marratech video-conferencing tool

The EDNA Data Model Framework

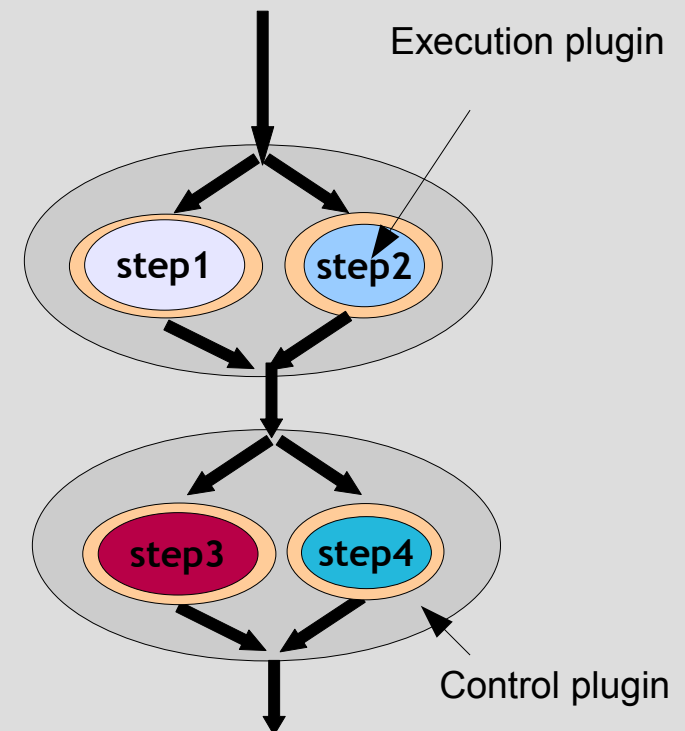
- From UML diagrams to generated code (data binding) :



```
<xs:element name="XSDDataSample" type="XSDDataSample">
  <xs:complexType name="XSDDataSample">
    <xs:complexContent>
      <xs:extension base="XSDData">
        <xs:sequence>
          <xs:element name="absorbedDose" type="XSDDataFloat" minOccurs="0" maxOccurs="1"/>
          <xs:element name="shape" type="XSDDataString" minOccurs="0" maxOccurs="1"/>
          <xs:element name="size" type="XSDDataFloat" minOccurs="0" maxOccurs="1"/>
          <xs:element name="susceptibility" type="XSDDataFloat" minOccurs="0" maxOccurs="1"/>
        </xs:sequence>
      </xs:extension>
    </xs:complexContent>
  </xs:complexType>
</xs:element>
```

EDNA Modularity : Plugins and their hierarchy

- Plugin base class :
 - Configuration, working directory, etc.
- Execution plugins :
 - Execution of external programs, e.g. (bash) scripts
- Controller plugins:
 - Control of execution plugins
 - Parallel execution
 - Synchronisation
- EDNA is based on AALib, however AALib is not a part of EDNA



EDNA Testing Framework

- The EDNA testing framework consist of three layers :
 - Kernel Unit tests
 - Plugin Unit tests
 - Plugin Execution tests

- Example of EDNA Plugin Execution tests result:

```
[SUCCESS] [ 1 ][ EDTestCasePluginExecuteReadImageHeaderv01.execute ][1.58020401001]
          [SUCCESS] [ 1 ][ EDTestCasePluginExecuteReadImageHeaderv01.testExecute ][1.41113901138]

[SUCCESS] [ 1 ][ EDTestCasePluginExecuteCCP4iv01.execute ][56.1572201252]
          [SUCCESS] [ 1 ][ EDTestCasePluginExecuteCCP4iv01.testExecute ][53.3898198605]
```

=====

```
[UnitTest]: #####
[UnitTest]: Summary Report:
[UnitTest]:               Total TestCases: 17
[UnitTest]:       Total TestCases [SUCCESS]: 17
[UnitTest]:       Total TestCases [FAIL]: 0
[UnitTest]:       [Total TestMethods]: 17
[UnitTest]:               Runtime: 341.2 [s]
[UnitTest]:               Run: 00d:00h:05m:41s:239ms
```

EDNA Collaborators

Alexander Popov^(d)

Alun Ashton^(e)

Andrew Leslie^(h)

Andrew McCarthy^(b)

Andrew Thompson^(k)

Clemens Schulze^(j)

Clemens Vornrhein^(f)

Darren Spruce^(d)

Elsbeth Gordon^(d)

Ezequiel Panepucci^(j)

Gérard Bricogne^(f)

Gerrit Langer^(b)

Gleb Bourenkov^(b)

Gordon Leonard^(d)

Harry Powell^(h)

Johan Turkenburg^(m)

Johan Unge^(g)

John Skinner⁽ⁱ⁾

Karl Levik^(e)

Katherine McAuley^(e)

Lucile Roussier^(k)

Marie-Françoise Incardona^(d)

Mark Basham^(e)

Meitian Wang^(j)

Michael Hellmig^(a)

Olga Roudenko^(k)

Peter Keller^(f)

Peter Turner^(l)

Pierre Legrand^(k)

Robert Sweet⁽ⁱ⁾

Romeu Pieritz^(d)

Sandor Brockhauser^(b)

Sean McSweeney^(d)

Takashi Tomizaki^(j)

Thomas Schneider^(b)

Uwe Mueller^(a)

(a) BESSY, Berlin, Germany

(b) EMBL, Grenoble, France

(c) EMBL, Hamburg, Germany

(d) ESRF, Grenoble, France

(e) Diamond Light Source, UK

(f) Global Phasing, Cambridge, UK

(g) MAX LAB, Lund, Sweden

(h) MRC LMB, Cambridge, UK

(i) NSLS, Brookhaven, U.S.

(j) SLS, Villigen, Switzerland

(k) Synchrotron Soleil, France

(l) University of Sydney, Australia

(m) University of York, UK

EDNA developers

Executive committee

Future developments

- Ranking in ISPyB
- Release of MXV1:
 - Based on the current prototype
 - Separate kernel
 - If time allows: implementation of data processing
- MXV2:
 - Improved generic data model
 - Implementation of XDS plugins
- Tomo:
 - EDNA for tomography

Why screen and rank?

- On average at the ESRF : 24 data collections per PDB deposition
- Hence it's important to collect data on the best crystals
- Screening : collect reference images from a large set of samples of the same structure and calculate a score for each sample
- Ranking : choose the sample with the highest score
- Without automation : screening is time-consuming and tedious
- Thanks to automation, screening and ranking are fast and easy

Screening and Ranking with EDNA

- Screening of an individual sample corresponds to a single DNA / EDNA characterisation
- Different philosophy of EDNA compared with DNA :
 - EDNA has no control of data collection, hence no GUI for screening
 - Screening will be handled by mxCuBE and other beamline control GUIs
- No implementation of ranking yet with EDNA
- Current idea : Ranking via the ISPyB interface
 - Advantages :
 - Ranking can be made independently of data collection
 - Same interface for DNA / EDNA
 - Disadvantage :
 - No connection (yet) of ranking results to the data collection software (mxCuBE)

Ranking via ISPyB

- Data collection page in the ISPyB GUI
- If implemented :
 - Several different screenings should be selectable and form the input for ranking

The screenshot displays the ISPyB (Information System for Protein Crystallography Baselines) web interface. The browser window is titled 'Data collections - Mozilla Firefox'. The URL is <http://ispyb.esrf.fr:8080/ispyb/ispyb/user/viewDataCollection.do?reqCode=displayForSession&sessionId=22653>. The page has a navigation bar with links: Lab-contacts, Shipment, Samples, Prepare experiment, Data collection, Feedback, and Help. The 'Data collection' link is active.

The main content area shows the 'Selected Session' details:

- Start Date: 06-11-2008
- BeamLine: ID14-1 33476882323
- Links: [Back to this session](#), [Back to sessions](#)

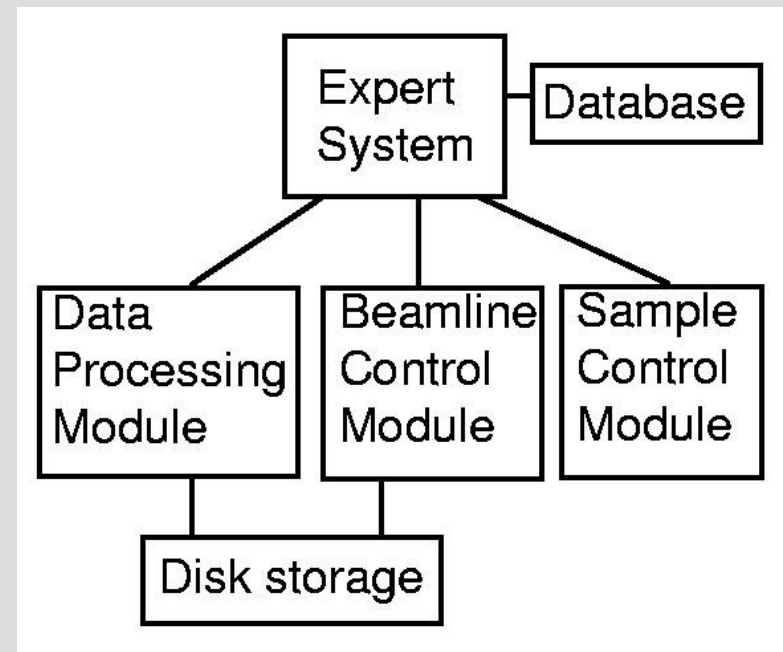
Below the session details is the 'Session Comments' section, which includes a text area for 'BeamLine Operator' and a 'Save' button. To the right of the comments section are links for 'View PDF report', 'PDF Screenings (2 images)', 'DOC Screenings (2 images)', and 'DOC Data Collections (>=3 images)'.

The 'Data Collections' table is shown below, with columns for various parameters. The table lists several data collection runs, including 'ref-mx415', 'mx415', '11-Xtal2', and 'postref'. A 'Save' button is visible next to the table.

Image Prefix	Run No	Start Time	# images	Wavelength Å	Transm.	Exposure Time sec.	Phi start °	Phi range °	Detector Resolution Å	Status	Snp/Run/Dna	Skip	Comments
ref-mx415	1	15:50:43	2	0.933	-1	1	0	1	2.48	●●●●			Collecting 2 reference images
mx415	4	15:11:17	10	0.933	-1	0.05	176	2.25	2.48	●●●●			DNA data collection
11-Xtal2	6	15:06:52	0			0.05	176	2.25	2.48	●●●●			DNA data collection
11-Xtal2	5	15:05:37	5	0.933	-1	0.05	176	2.25	2.48	●●●●			DNA data collection
postref	5	15:05:10	3	0.933	-1	0.05	266	2.25	2.48	●●●●			

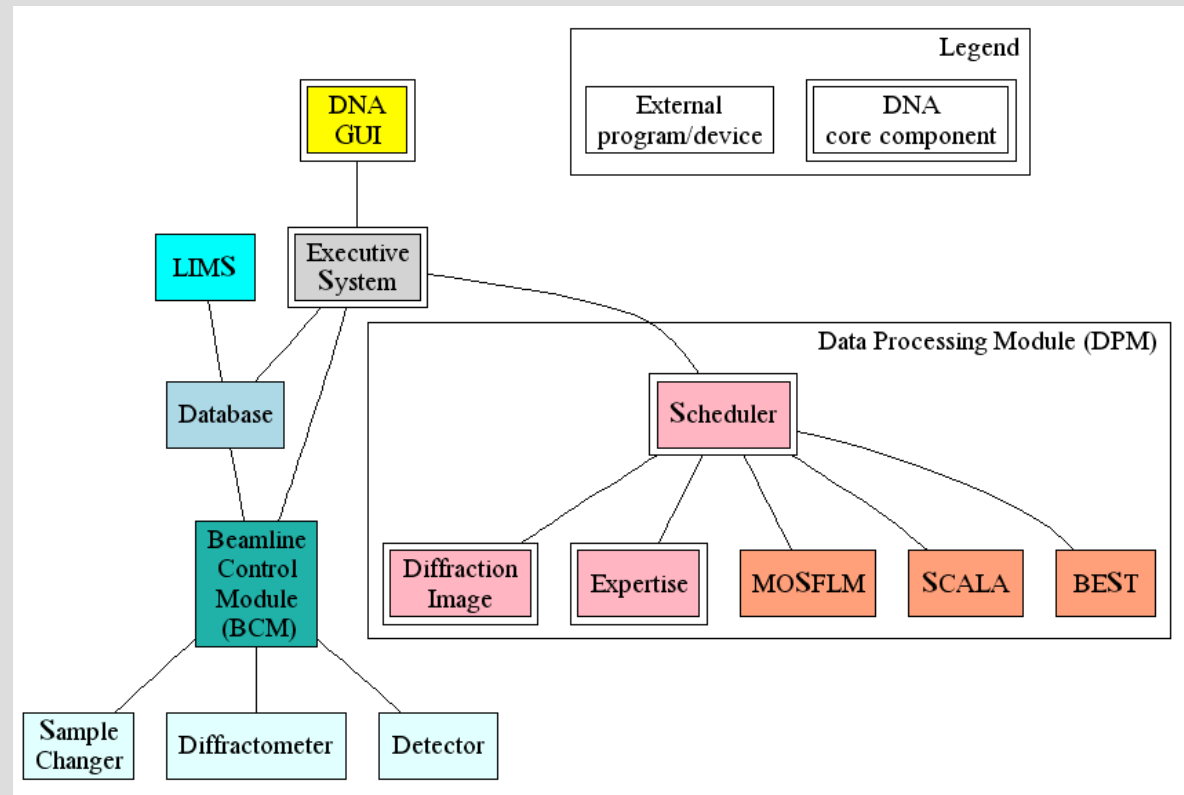
The DNA collaboration - the beginning

- Kick-off meeting in 2001
- Initial collaborators :
 - ESRF
 - Daresbury SRS
 - MRC LMB Cambridge
- Initially no external funding
- Meaning of "DNA" :
 - automated collection of data



Evolution of the DNA collaboration

- Main development period 2001 - 2005
- More collaborators and more developers entered the project, mainly thanks to external fundings : BioXHIT and e-HTPX.
- Installed and used at :
 - ESRF
 - Diamond
 - Recently used at NSLS (Brookhaven), now using EDNA
- Part of the 2008 BESSY Innovation Award



DNA collaborators

DNA Collaborators (in 2007)

Home institute	Name and link to email	DNA related work funded by	DNA tasks
Brookhaven Nat'l Lab.	Alex Soares	BNL	DNA co-ordinator at BNL
	Bob Sweet		
	John Skinner		
Diamond Light Source	Alun Ashton	DLS	Acting DNA project co-ordinator
	Colin Nave		
	Elizabeth Duke		DNA co-ordinator at DLS
	Karl Levik		DNA developer
	Katherine McAuley		Testing and implementation of DNA at Diamond Light Source
EMBL Grenoble	Raimond Ravelli	EMBL Grenoble	DNA co-ordinator at the EMBL Grenoble
	Sandor Brockhauser	BioXHIT	Working on introducing kappa geometry strategy into DNA
EMBL Hamburg	Alexander Popov	BioXHIT	DNA co-ordinator at EMBL Hamburg
	Gleb Bourenkov	DESY	
	Venkataraman Parthasarathy	SPINE	Integration of the BEST strategy software
	Sean McSweeney	ESRF	DNA co-ordinator at the ESRF
	Darren Spruce		Responsible for the ESRF BCM (ProDC) and the DNA - LIMS connection
ESRF	Olof Svensson		Responsible for the DNA Executive System and DNA 2.0 Project Manager
	Marie Francoise Incardona	BioXHIT	Working on DNA 2.0
	Romeu Pieritz		Responsible for developing the ranking module and working on DNA 2.0
Global Phasing	Gérard Bricogne	Global Phasing	DNA co-ordinator at Global Phasing
	Peter Keller	BioXHIT	Working on DNA 2.0
MRC LMB Cambridge	Andrew Leslie	MRC	DNA co-ordinator at Cambridge
	Harry Powell	CCP4	Responsible for the DNA DPM based on MOSFLM
MRC France - BM14 at the ESRF	Ludovic Launier	e-htpx	Working together with Darren on the DNA - beamline database connection
STFC Daresbury	Graeme Winter	e-htpx	Responsible for the DNA Scheduler
SLS - PSI	Takashi Tomizaki	SLS - PSI	DNA co-ordinator at SLS
Synchrotron Soleil	Andrew Thompson	Synchrotron Soleil	DNA co-ordinator at Soleil
	Lucile Roussier		Working on DNA - database connection for Soleil
	Eric Girard		Working on offline tests of DNA
	Pierre Legrand		Working on integrating XDS as a DPM in the DNA system

What we got right in DNA

- Collaboration :
 - Scientists involved in design and testing
 - Executive committee for setting milestones / deliverables and for resolving conflicts between developers
- Two major Use Cases implemented :
 - Characterisation + data collection with online integration and quick scaling
 - Automatic screening and ranking
- Used regularly at DLS, ESRF and recently at the NLSL (now replaced by EDNA)

What went wrong in DNA

- The choice of name...
- Collaboration :
 - No project agreement
 - Minimal project management
- Not modular:
 - Too costly to change work flow
 - Poorly designed data model
 - Difficult for new developers to enter the collaboration
- MX hardwired!

A new Project!

Initial Design Thoughts

- A framework for Online Data Analysis
 - Use-case driven development
 - Modular - based on plugins
 - Configuration facility
 - Testing framework for assuring robustness
 - Data model tools and data classes code generation from UML
- Project management:
 - Executive committee, Project coordinator, Project manager
 - Code style / code reviews
- Application to other scientific online data analysis tasks
 - No MX specific code in the kernel

DNA 2.0 → EDNA

- The new project was accepted by the DNA executive committee Autumn 2005
- Project Manager
- DNA 2.0 officially launched in February 2007 in a DNA meeting held here at the DLS
- New name EDNA decided in the Project Agreement meeting at the ESRF in October 2007

EDNA Developments and Events

- February 2007 : Launch of DNA 2.0
- June 2007 : The spike meeting (ESRF)
- October 2007 : The Project Agreement meeting (ESRF)
- November 2007 : The first EDNA developers' workshop
- February 2008 : Developers' meeting at Soleil, France
- June 2008 : The prototype demonstration meeting (ESRF)
- August 2008 : Release of the EDNA prototype
- October 2008 : Plugin developers' workshop (ESRF)

Conclusions

- EDNA is a framework for online data analysis
 - Modular
 - Data modelling framework
 - Testing framework
- EDNA is a collaboration
 - Project management
 - Many facilities/institutes participating in the collaboration
- The EDNA prototype is ready to use
- The EDNA MXV1 application will be released in May
- MXV2 / Tomography applications are being planned

Acknowledgements

- The DNA team
- The EDNA team
- The ISPyB team
- Sasha Popov for providing me his slides
- DLS for inviting me and hosting this EDNA developers' meeting
- ...and thank you for your attention!