# EDNA : A framework for Advanced On-line Data Analysis

Olof Svensson

Data Analysis Unit / Instrument Support and Development Division

ESRF

# Outline

- Advanced On-line Data Analysis
- EDNA
    - The collaboration
    - The framework
- Current ODA applications based on EDNA
    - MX
    - non-MX
- Future developments
    - Workflow tools

# On-line Data Analysis I

- Data analysis performed as quickly as possible in order to provide feedback to the user(s):
  - Quick calculations for data visualisation
  - Quality of the data generated by the experiment – is the experiment successful?
  - Use of results for planning the experiment : e.g. MX strategy calculation

- Data analysis during experiment
  - Users leave with reduced data (in addition to raw data)

# On-line Data Analysis II

- Combination of data acquisition and data processing

- Automation
  - Software - and hardware
  - Robustness crucial

- Impact on beamline efficiency :
  - Should introduce added value
  - No slow-down of experiments
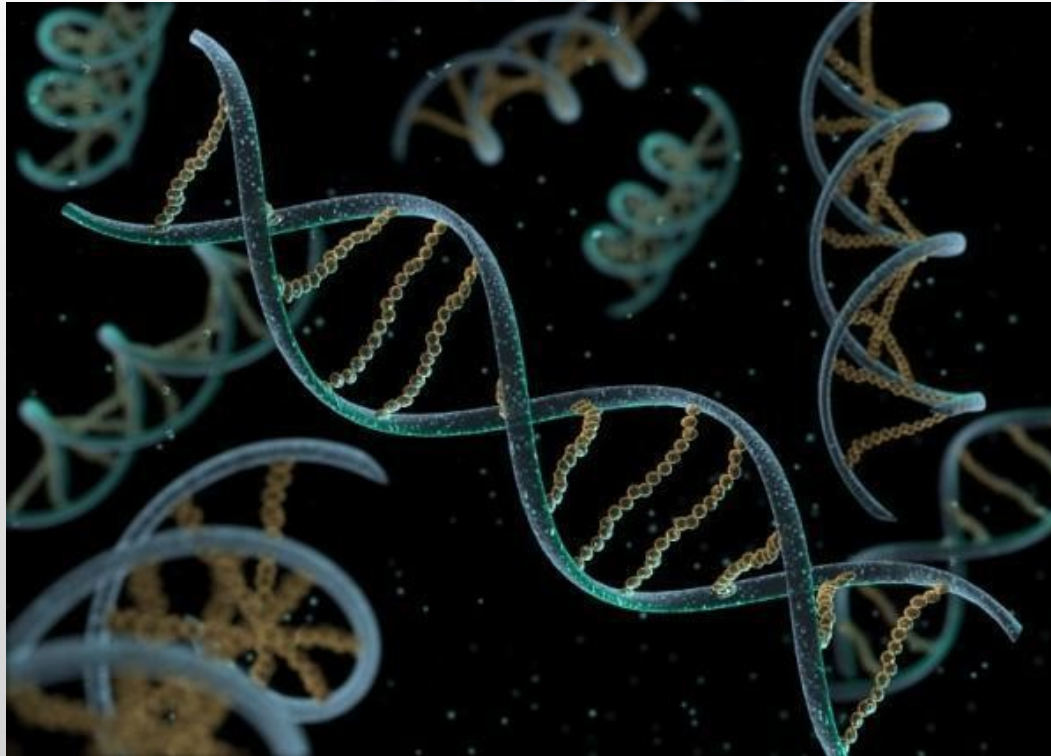  - No or little impact on beamline scientists / users workload

# Advanced On-line Data Analysis - Examples

- Non-trivial data analysis :
  - Azimuthal integration followed by peak fitting and/or integration
  - Tomographic reconstruction

- Complex feedback to data acquisition :
  - MX data collections strategy
  - Modification of data collection strategy during data acquisition

- A combination of the two :
  - MX strategy data collections taking into account radiation damage and kappa geometries

# On-line data analysis: Challenges

- Moving target:
  - What's today's state-of-the-art science might not be so tomorrow

- Limited resources:
  - Re-use existing scientific software
  - Collaborate
  - → Avoid re-inventing the wheel

- Automation
  - Robustness crucial
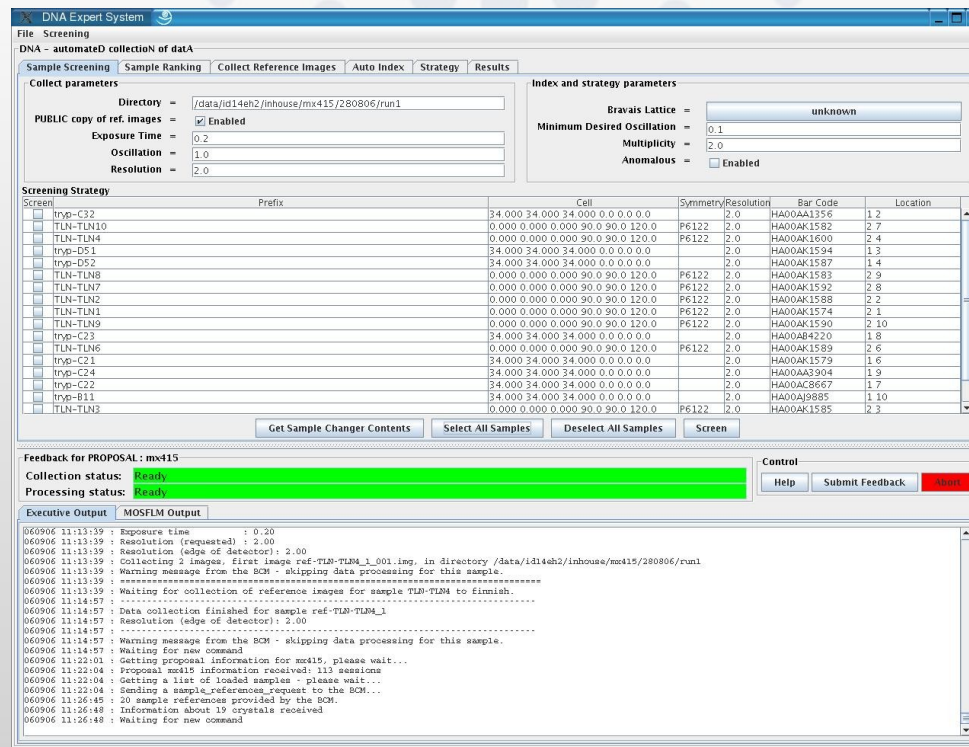
- Conflicting challenges!

# In the beginning there was DNA...



… or RNA or something else, if you speak about life...

# In the beginning there was DNA...



## … if you speak about EDNA!

# Beginning of automated MX characterisation...
## … or DNA in one slide

- Kick-off meeting in 2001
- Initial collaborators :
  - ESRF
  - Daresbury SRS
  - MRC LMB Cambridge
- Initially no external funding
- Meaning of "DNA" :
  - automate**D** collectio**N** of dat**A**
- **Main development period 2001 – 2005**
- More collaborators and more developers entered the project, mainly thanks to external fundings : BioXHIT and e-HTPX.
- **Major component of the 2008 BESSY innovation award**

# Why DNA became EDNA

- Many positive experiences learned from the DNA project :
    - Collaborative project : scientists / developers from several facilities working towards a common goal
    - Milestones achieved : automated MX characterisation

- However also many negative experiences :
    - The choice of name...
    - No project agreement, minimal project management
    - Not modular : difficult to add new features, difficult for new developers to enter the project
    - MX hardwired

# EDNA

- Project started in 2007

- Many collaborators and ideas brought over from the DNA project however no shared code

- Major development goals - collaborative developments :
  - Project agreement, project management
  - Workflow developments
  - Modular
  - Data model framework
  - Testing framework
  - Not hard-wired to MX

# EDNA Collaboration

- EDNA is about collaboration:
  - Code sharing (SVN)
  - Coding conventions
  - Code reviews
  - Open source (LGPL, GPL)
  - Bug tracker
  - Wiki : http://www.edna-site.org
  - Memorandum of Understanding
  - Executive committee
  - Project manager / coordinator
  - Regular meetings / video conferences

European Synchrotron Radiation Facility

# 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

Execution plugin

step1    step2

step3    step4

Control plugin

# Example EDNA workflow : MXv1 Characterisation

- MX sample characterisation taking into account radiation damage

- Indexing using MOSFLM or Labelit

- Parallel integration of reference images

- If flux + beamsize:
  - RADDOSE for estimating radiation damage
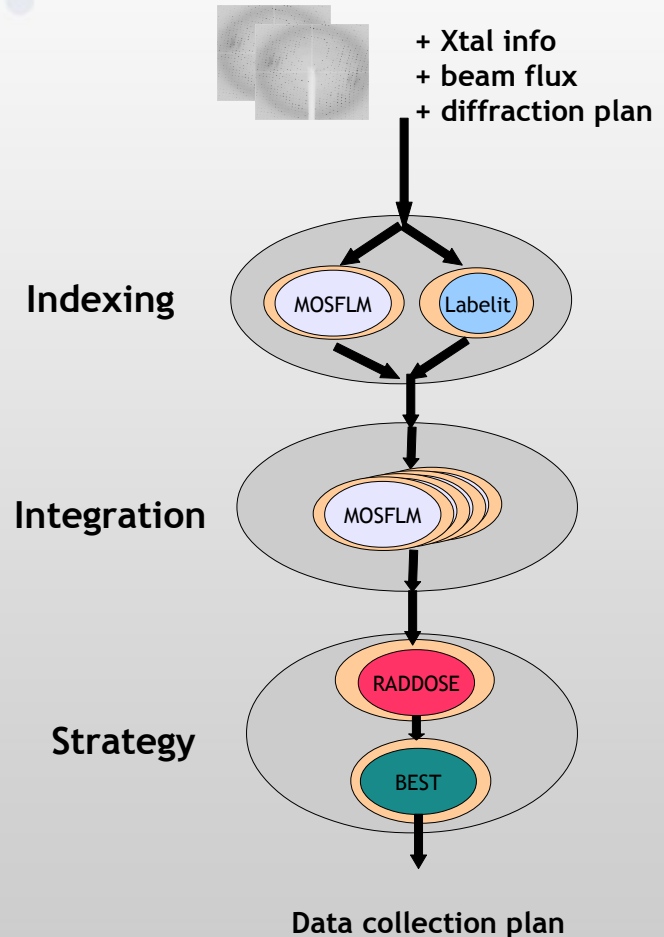
- BEST strategy calculation
  - taking into account radiation damage
  - multi-subwedge data collection strategies



+ Xtal info
+ beam flux
+ diffraction plan

**Indexing** — MOSFLM, Labelit

**Integration** — MOSFLM

**Strategy** — RADDOSE, BEST

Data collection plan

# EDNA Data Modelling Framework

- Experiments produce data
- Data cannot in general be analysed without meta-data
- Complex experiments → complex meta-data

# EDNA Testing Framework

- Unit, execution and regression tests:

```
[UnitTest]: #######################################################################
[UnitTest]: Result for EDTestSuitePluginExecuteAll : FAILURE
[UnitTest]:
[UnitTest]:  Number of executed test suites in this test suite : 5
[UnitTest]:
[UnitTest]:
[UnitTest]:   Total number of test cases executed with SUCCESS : 124
[UnitTest]:   Total number of test cases executed with FAILURE : 8
[UnitTest]:
[UnitTest]:
[UnitTest]: OBS! The following test methods ended with failure:
[UnitTest]:
[UnitTest]:   EDTestCasePluginExecuteControlCharForReorientationv2_0_noKAPPA_ :
[UnitTest]:    testExecute :
[UnitTest]:       Plugin failure assert: should be False, was True FAILURE: Expected different from obtained - identifier
/mntdirect/_scisoft/users/svensson/tmp/EDTestSuitePluginExecuteAll_20110114-093729/tmpPW7olu
[UnitTest]:
 ...
[UnitTest]:
[UnitTest]:
[UnitTest]: Total number of test methods executed with SUCCESS : 129
[UnitTest]: Total number of test methods executed with FAILURE : 8
[UnitTest]:
[UnitTest]:                            Runtime : 1108.997 [s]
[UnitTest]: #######################################################################
```
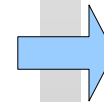
- Test suites executed automatically every night

- Continuous integration: nightly builds

# Zen of EDNA

- Robustness is more important than performance
- Modularity brings maintainability
- Data-structures should be separated from code
- Flexibility comes with modularity & data-structures
- Community and collaboration to fight bugs
- Data parallelism above multi threading
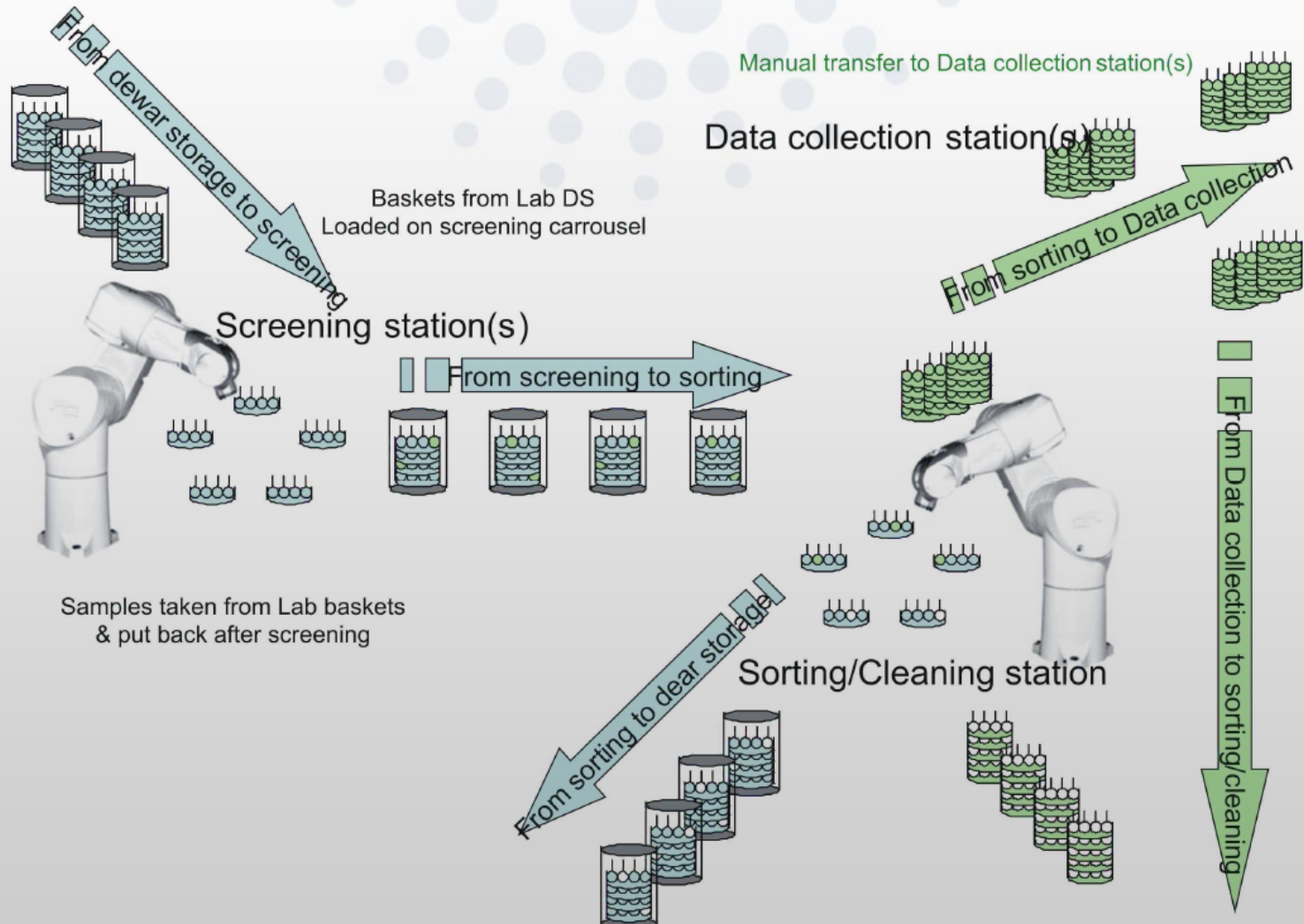- Performance & parallelism are provided by the EDNA-kernel

# Edna toolbox

- > 100 execution plugins:
  - Generic command line execution, Image conversion, movied ...
  - HDF5 writers for stack of images, map of spectra
  - Conditional branching, accumulator of information
- 5 real applications available from repository (2 demo)
  - BioSaxs
  - DiffractionCT
  - Dimple
  - Fullfield XANES
  - MX v1 & v2
  - Demo projects: Ccp4 and Raw photography
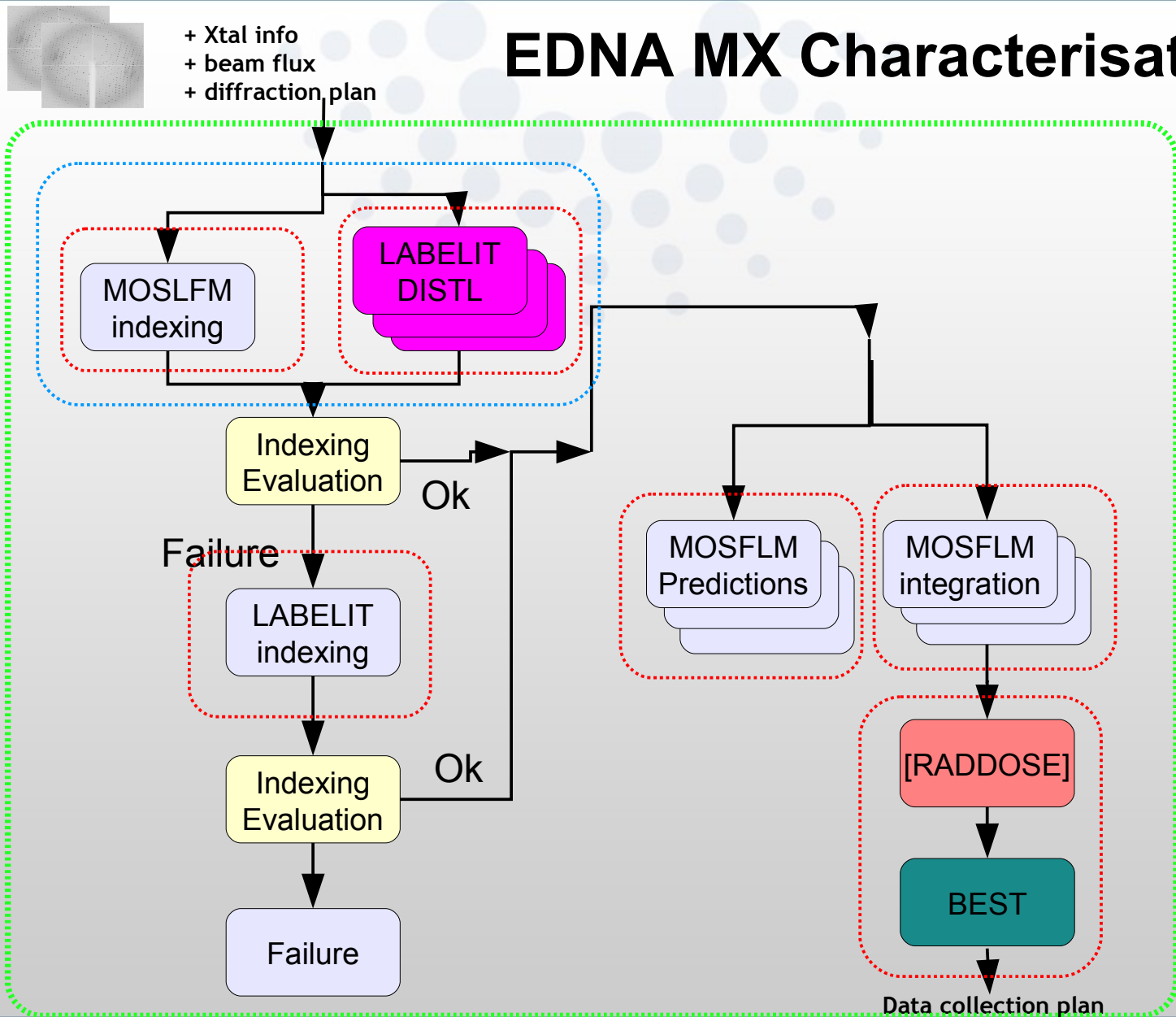- 3 kind of launchers: Command line, Parallel, Tango

# Scientific EDNA workflows @ ESRF

- Macromolecular crystallography:
  - Characterisation taking into account radiation damage (MOSFLM, Labelit, RADDOSE, BEST)
  - Connection with experiment data base (ISPyB)
  - Parallel execution of characterisation (GRID data processing)
  - Parallel creation of image thumbnails
- Diffraction Computed Tomography
  - SPD: Image correction, fast azimuthal integration
  - Sinograms saved in HDF5 format
- Small Angle Scattering
  - Image correction and fast azimuthal integration
- Full Field XAS
  - Image correction (dark, flat)
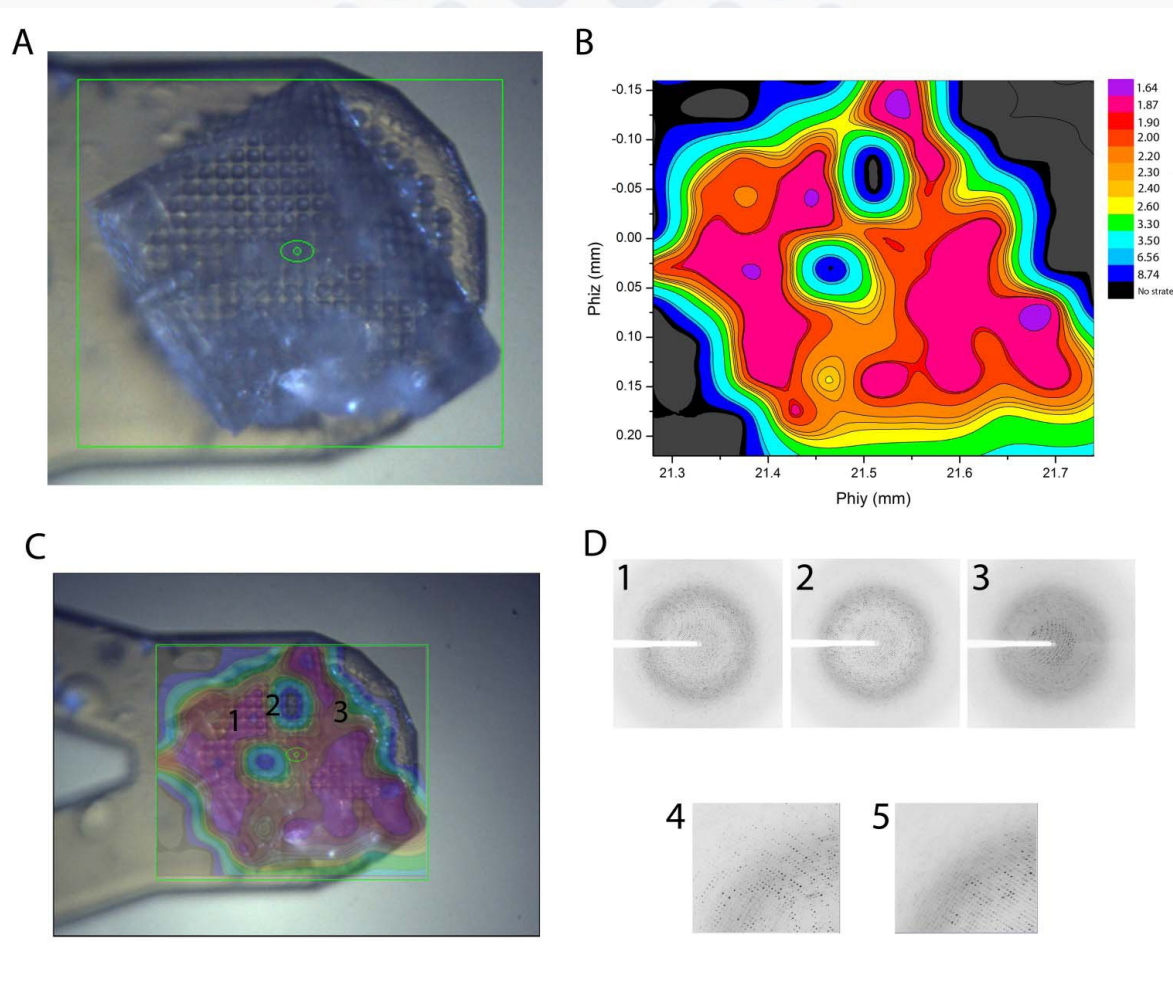  - Image alignment (offset measurements by FFT)
  - HDF5 output

# Challenge for the ESRF Upgrade :
# Massively Automated Sample Selection Integrated Facility
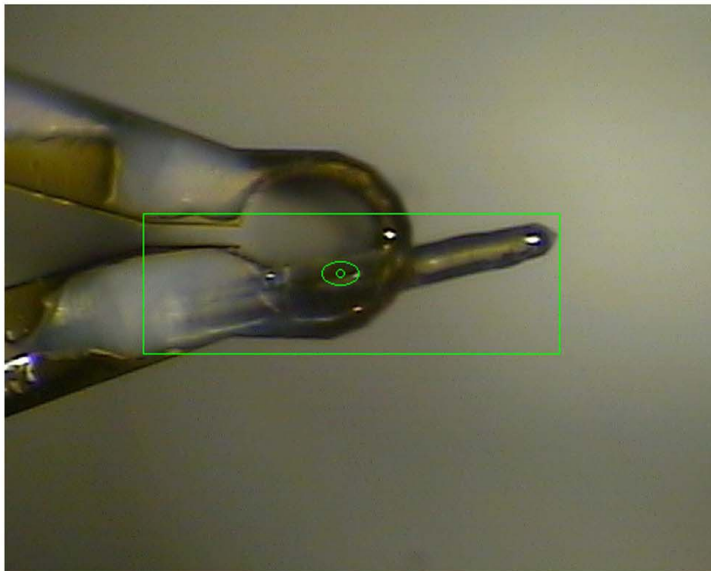
# EDNA MX Characterisation



+ Xtal info
+ beam flux
+ diffraction plan

MOSLFM indexing

LABELIT DISTL

Indexing Evaluation

Ok

Failure

LABELIT indexing

Indexing Evaluation

Ok

Failure

MOSFLM Predictions

MOSFLM integration

[RADDOSE]

BEST

Data collection plan

# MX Grid Scans I

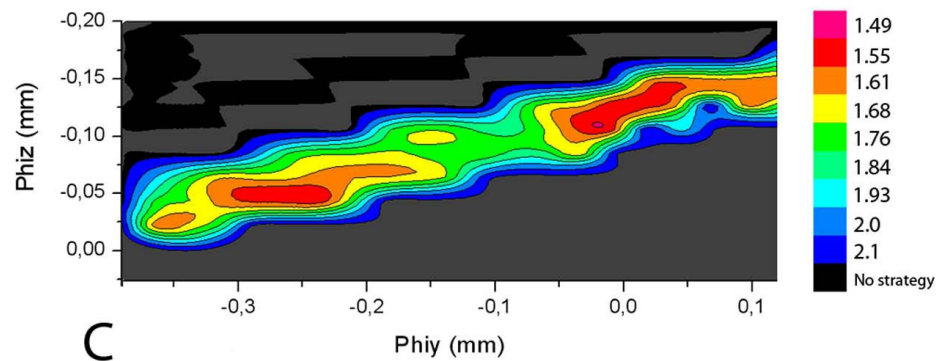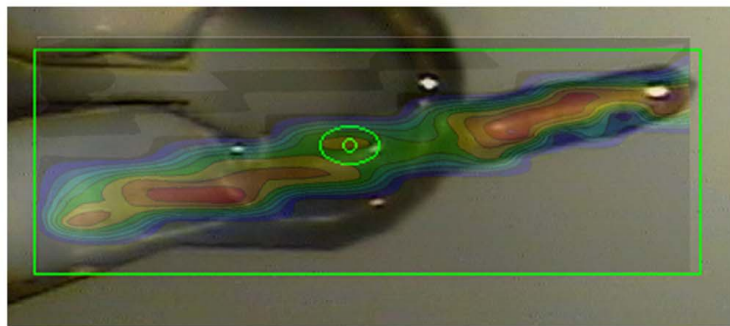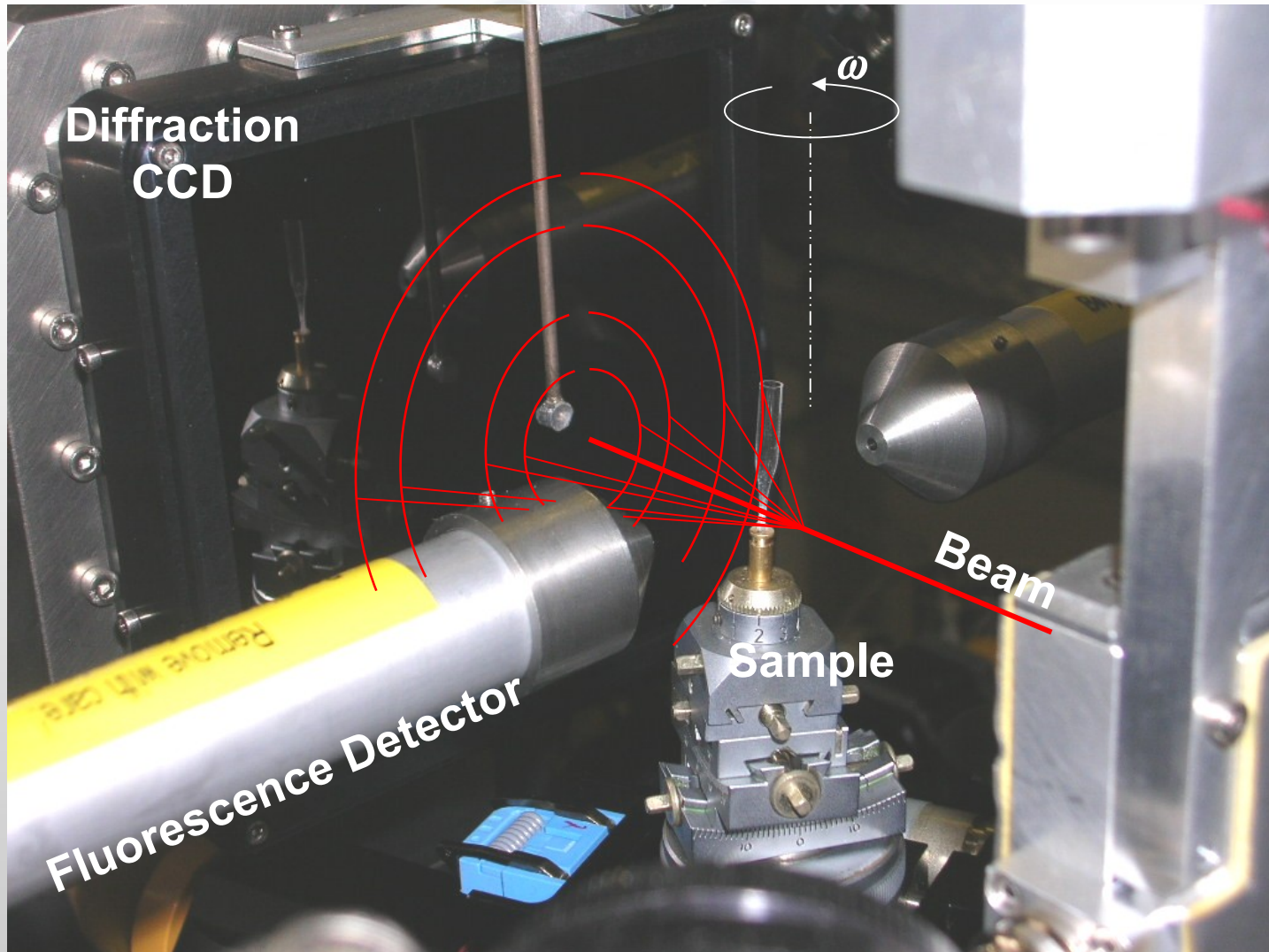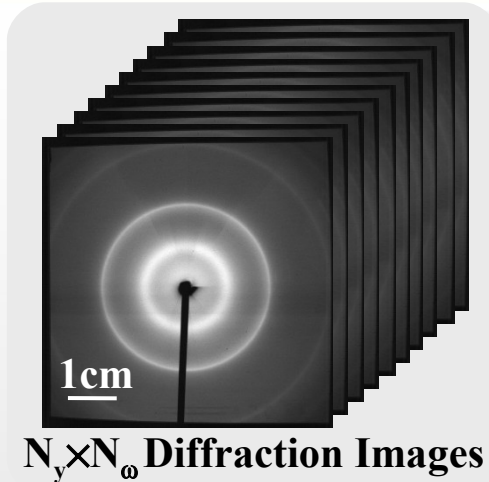# MX Grid Scans II



Bowler M.W. , Guijarro M. , Petitdemange S. , Baker I. , Svensson O. , Burghammer M. , Mueller-Dieckmann C. , Gordon E.J. , Flot D. , McSweeney S.M. , Leonard G.A. - Diffraction cartography: Applying microbeams to macromolecular crystallography sample evaluation and data collection Acta Crystallographica D 66, 855-864 (2010)
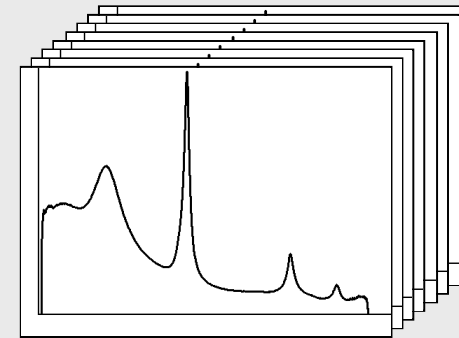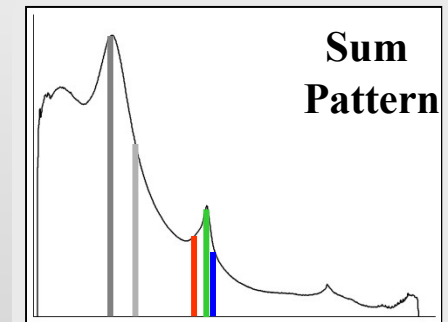
## ID22 – Fluorescence-Diffraction Tomography



Diffraction CCD

$\omega$

Beam
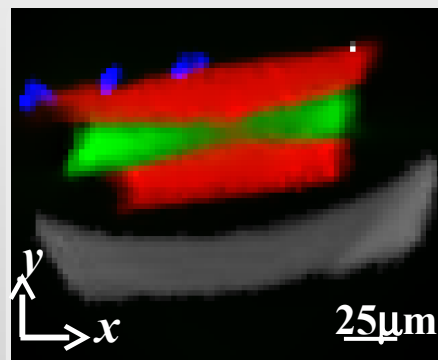
Sample

Fluorescence Detector

$z$

$x$

$y$

1cm

Courtesy of V.A. Solé and J. Kieffer

**Azimuthal Integrations**

**Fit2d software**

$N_y \times N_\omega$ **Diffraction Images**

1cm

$N_y \times N_\omega$ **Diffraction Patterns**

**Sum Pattern**

$y$

$\omega$

**Phase Sinograms**

**PyMca software**

**Sum Sinogram**

$\omega$

**Reconstruction**

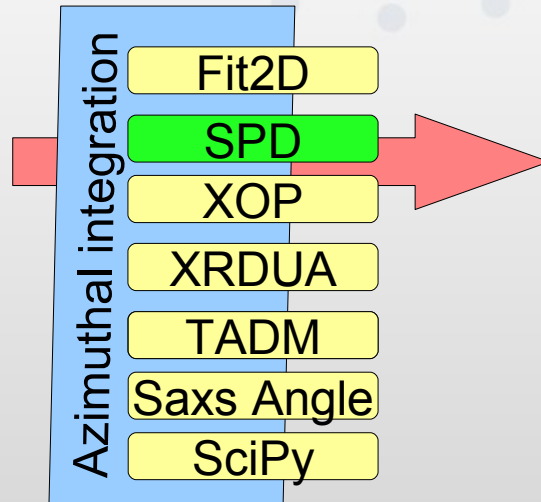| | | | |
|---|---|---|---|
| ⬜ (gray) | *Capillary* | 🟦 | *Ferrite* |
| ⬜ (white) | *Calcite* | 🟥 | *sp3* |
| | | 🟩 | *Cubic* |

$y$

$x$

25μm

Acknowledgements: Pierre Bleuet CEA - Grenoble

# Online data Analysis for NINA

- Reduce each image (2D) to a spectrum (1D)
- Store a 2D mapping of spectra

**N** x **ω** 1D diffraction patterns



**N** x **ω** 2D
Diffraction
Images

Azimuthal integration: Fit2D, SPD, XOP, XRDUA, TADM, Saxs Angle, SciPy
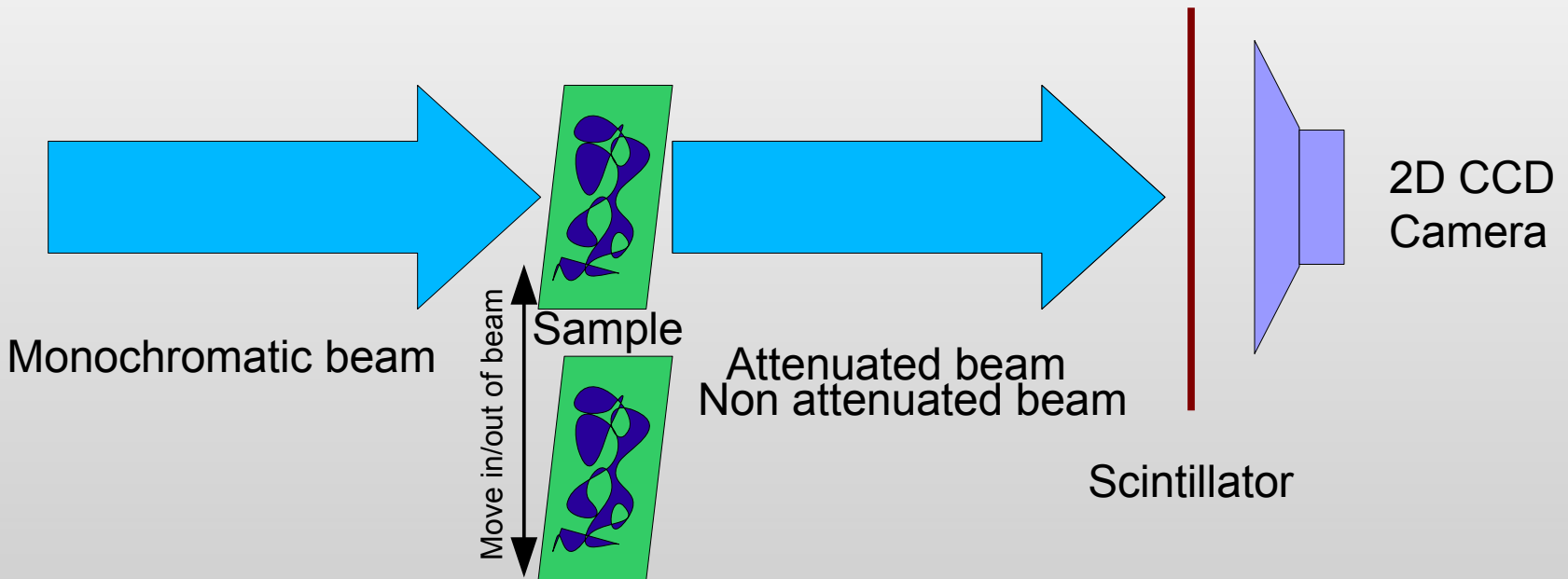
Diffraction 2θ

ω Angles

N points

Performance: 4 images / s, full treatment in 40 min.

Courtesy of J. Kieffer

# FullField XAS

Scan in Energy on monochromator
around absorption edge of a given element

Sample size: 1 mm x 1 mm x 1 μm
Resolution: 500 nm



Monochromatic beam

Move in/out of beam

Sample

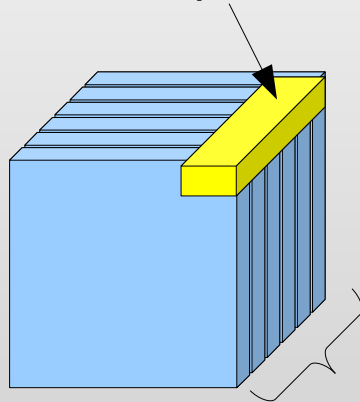Attenuated beam
Non attenuated beam

Scintillator

2D CCD
Camera

→ Measure a (couple of) flats at each energy to correct for scintillator's response
→ Align the sample at each position to correct sub-micron position change
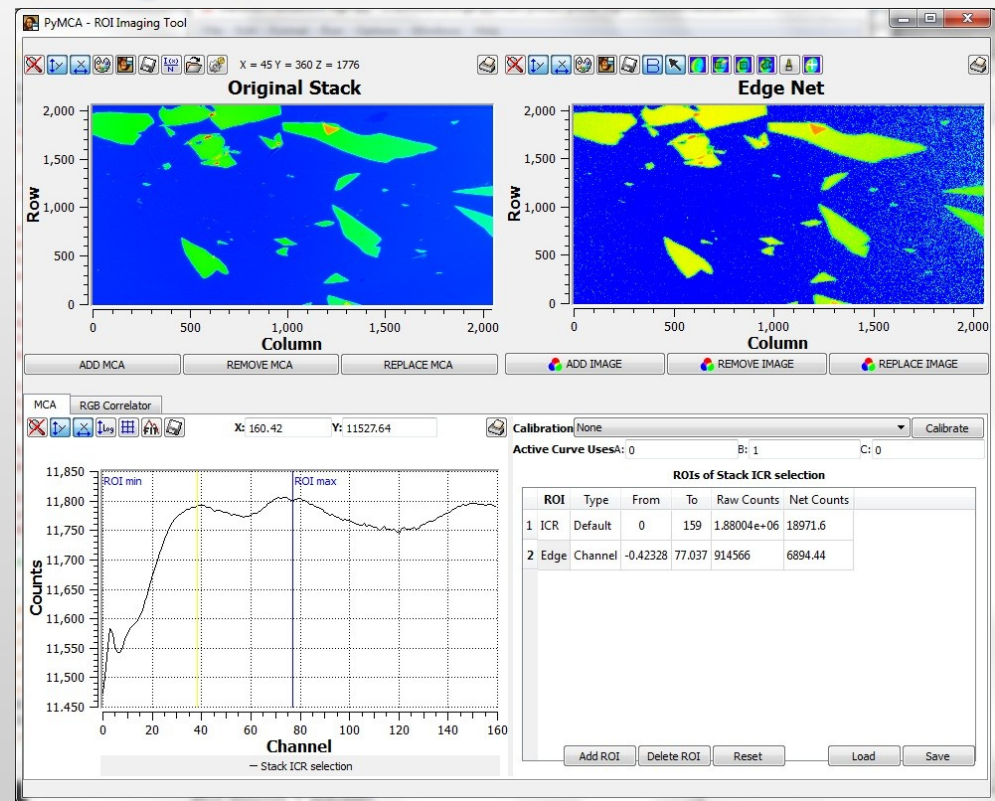
Courtesy of J. Kieffer

# FullField XAS Results

- Tested online & offline in December 2010 (ID21)
- Image Alignment by FFT OK ($\to$ ID10)
- Nexus compliance

Xanes Spectra

Stack of images



Courtesy of J. Kieffer and A. Sole

# Workflow tools



DOI:10.1145/1897852.1897871

**Compose "dream tools" from continuously evolving bundles of software to make sense of complex scientific data sets.**

BY KATY BÖRNER

# Plug-and-Play Macroscopes

DECISION MAKING IN science, industry, and politics, as well as in daily life, requires that we make sense of data sets representing the structure and dynamics of complex systems. Analysis, navigation, and management of these continuously evolving data sets require a new kind of data-analysis and visualization tool we call a macroscope (from the Greek macros, or "great," and skopein, or "to observe") inspired by de Rosnay's futurist science writings.[8]

Just as the microscope made it possible for the naked human eye to see cells, microbes, and viruses, thereby advancing biology and medicine, and just as the telescope opened the human mind to the immensity of the cosmos and the conquest of space—the macroscope promises to help make sense of yet another dimension—the infinitely complex. Macroscopes provide a "vision of the whole," helping us "synthesize" the related elements and detect patterns, trends, and outliers while granting access to myriad details.[18,19] Rather than make things larger or smaller, macroscopes let us observe what is at once too great, slow, or complex for the human eye and mind to notice and comprehend.

Many of the best micro-, tele-, and macroscopes are designed by scientists keen to observe and comprehend what no one has seen or understood before. Galileo Galilei (1564–1642) recognized the potential of a spyglass for the study of the heavens, ground and polished his own lenses, and used the improved optical instruments to make discoveries like the moons of Jupiter, providing quantitative evidence for the Copernican theory. Today, scientists repurpose, extend, and invent new hardware and software to create macroscopes that may solve both local and global challenges[20] (see the sidebar "Changing Scientific Landscape").

My aim here is to inspire computer scientists to implement software frameworks that empower domain scientists to assemble their own continuously evolving macroscopes, adding and upgrading existing (and removing obsolete) plug-ins to arrive at a set that is truly relevant for their work—with little or no help from computer scientists. Some macroscopes may resemble cyberinfrastructures (CIs),[1] providing user-friendly access to massive amounts of data, services, computing resources, and expert communities. Others may be Web services or stand-alone tools. While microscopes and telescopes are physical instruments, macroscopes resemble continuously changing bundles of software plug-ins. Macroscopes make it easy to select and combine algorithm and tool plug-ins but also interface plug-ins, workflow support, logging, scheduling, and other plug-ins needed for scientifically rigorous work. They make it easy to share
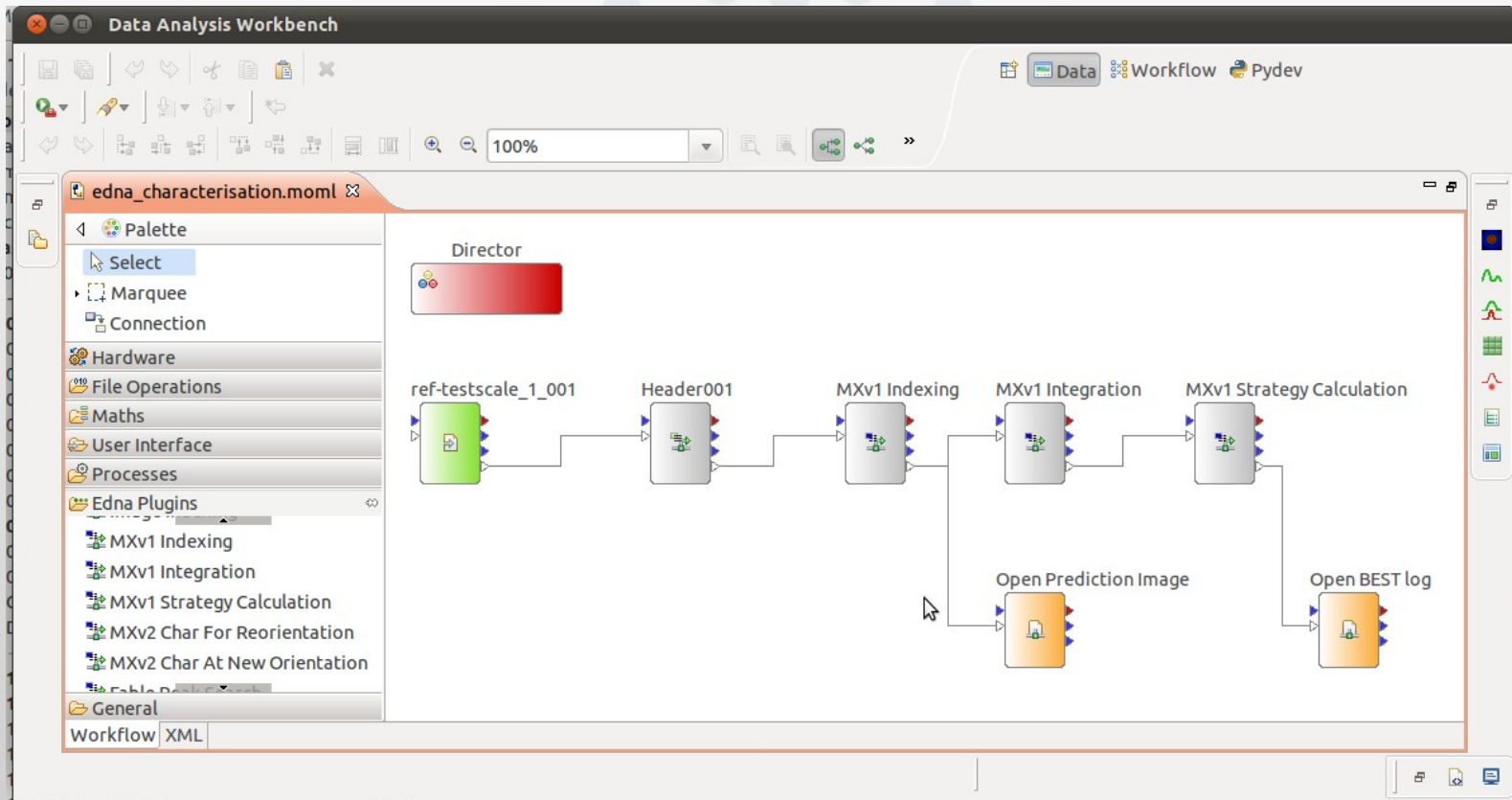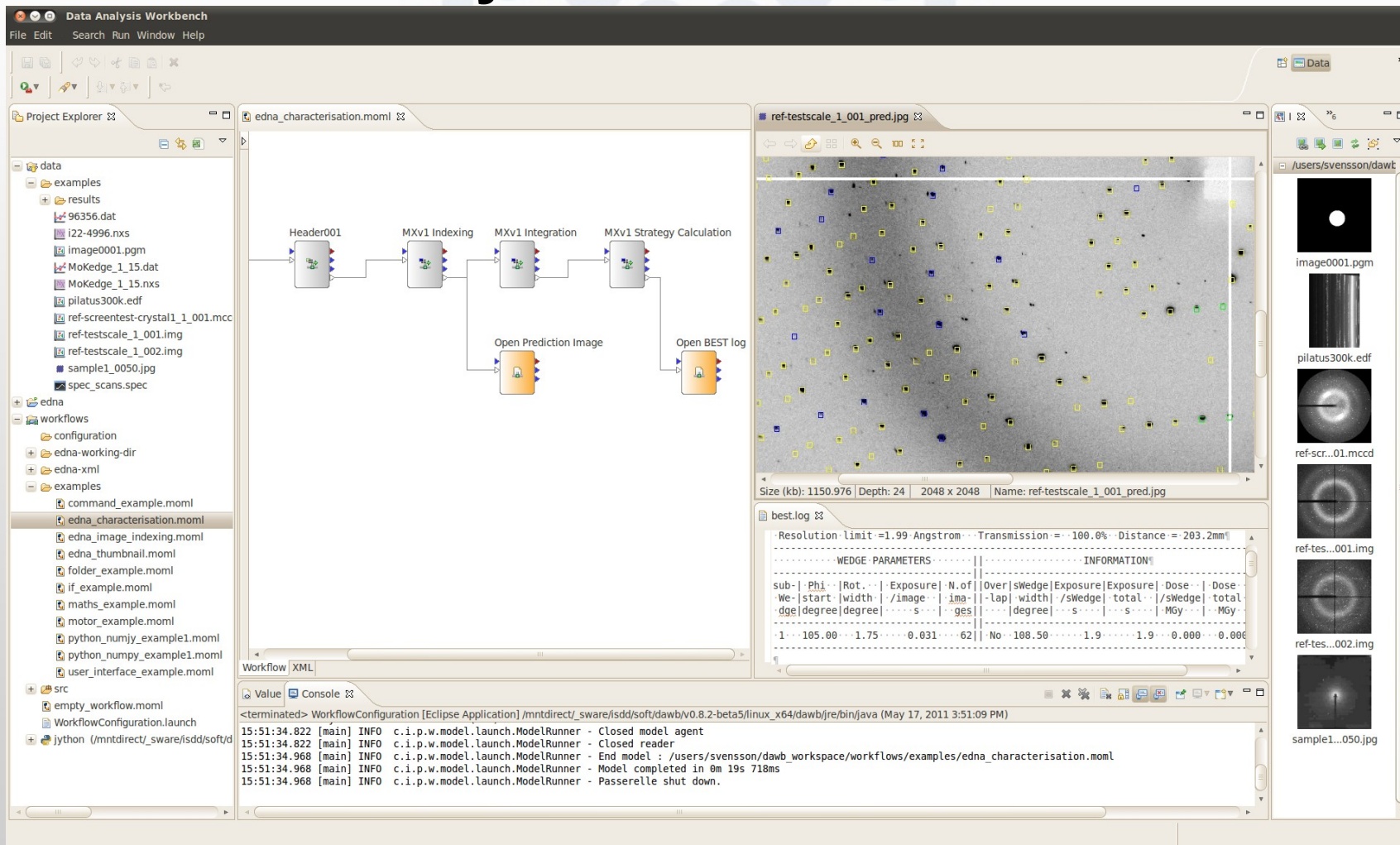
» key insights
- OSGi/CIShell-powered tools improve decision making in e-science, government, industry, and education.
- Non-programmers can use OSGi/CIShell to assemble custom "dream tools."
- New plug-ins are retrieved automatically via OSGi update services or shared via email and added manually; they can be plugged and played dynamically, without restarting the tool.

My aim here is to inspire computer scientists to implement software frameworks that empower domain scientists to assemble their own continuously evolving macroscopes, adding and upgrading existing (and removing obsolete) plug-ins to arrive at a set that is truly relevant for their work—with little or no help from computer scientists.

http://cacm.acm.org/magazines/2011/3/105316-plug-and-play-macroscopes/fulltext

# EDNA Workflow Tool - Example

# Data Analysis WorkBench - DAWB



http://www.dawb.org

# Acknowledgements