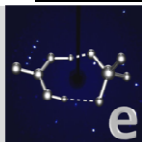# The eCrystals Federation

**Repository Curation Service Environments (RECURSE) Workshop**
**National e-Science Centre, Edinburgh**

**4th International Digital Curation Conference**
**"Radical Sharing: Transforming Science?"**
**1-3rd December 2008**
**Edinburgh, Scotland**

**Manjula Patel**
**UKOLN, University of Bath, UK**

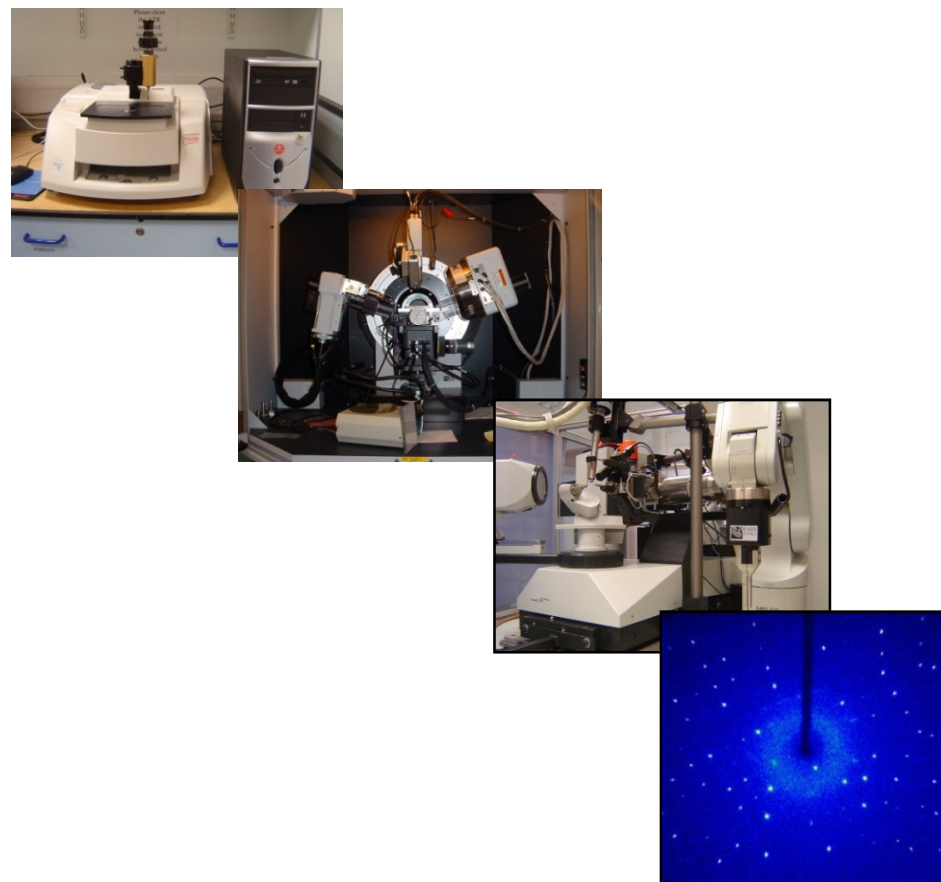eCrystals Federation

# Context

- The data deluge
  - Advances in instrumentation, data storage technologies, computational power and improvements in algorithms
  - Development of grid and cyber infrastructures
- Actual nature of science is changing
  - Mining and analysis of large datasets (e.g. Protein Data Bank, GenBank)
  - Open Science (e.g.Open Notebook Science; myExperiment)
- High quality data are the raw materials of contemporary e-science
  - Verification; Validation; Replication
  - Predictive science
  - Innovative scientific endeavour
- S. Carlson, *Lost in a Sea of Science Data*, The Chronicle of Higher Education, June 2006

  "To vet experiments, correct errors, or find new breakthroughs, scientists desperately need better ways to store and retrieve research data"

  "Data from Big Science is … easier to handle, understand and archive. Small Science is horribly heterogeneous and far more vast. In time Small Science will generate 2-3 times more data than Big Science."
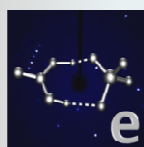
# Crystallography –The Science

- Sub-discipline of chemistry
- Concerned with determining the structure of a molecule and its 3D orientation with respect to other molecules in a crystal
- Analysis of diffraction patterns obtained from X-ray scattering experiments
- Focus on laboratory based experimental technique of chemical crystallography undertaken at the EPSRC National Crystallography Service (NCS), UK
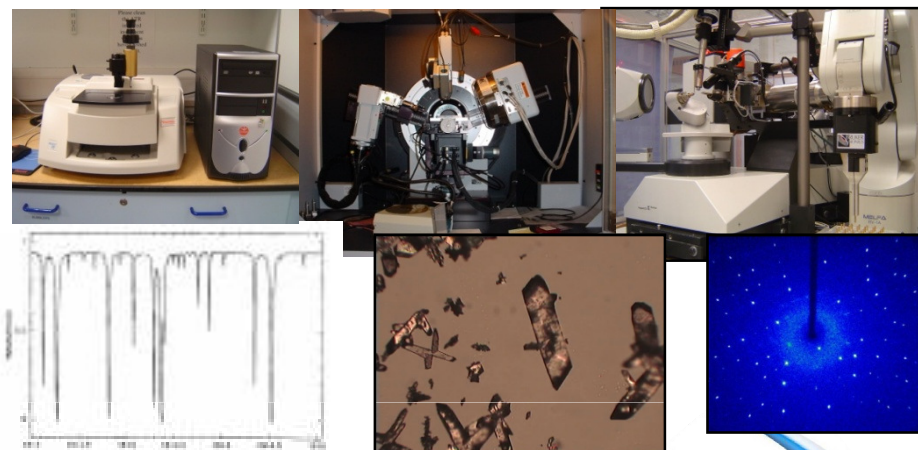


Images from Simon Coles (NCS), 2006
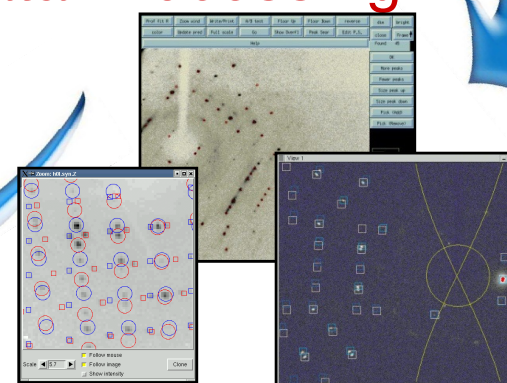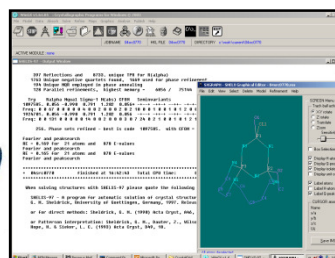
**eCrystals** Federation

# Data Generation

**Synthesis**

**Data Collection**

**Publication**

**Data Processing**

**Data Workup**

Cambridge Crystallographic Data Centre
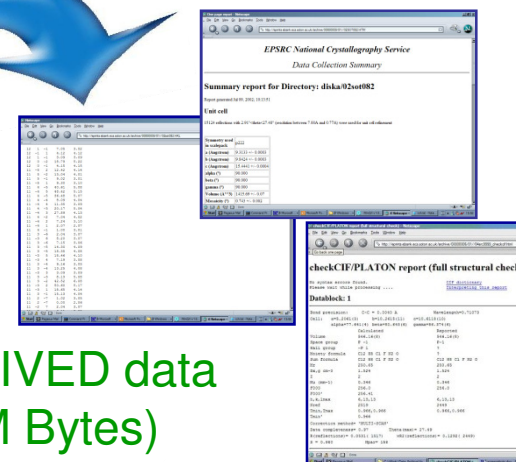
Adapted from Simon Coles (NCS), 2007

**eCrystals** Federation

# Data Volumes



RAW data
(G Bytes)

DERIVED data
(M Bytes)

Laboratory; Institution

Subject Repository; Data Centre; Public Domain

RESULTS data
(K Bytes)

**eCrystals** Federation

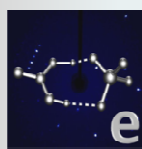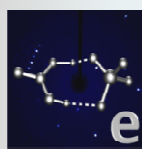# Community & Current Practice (1)

- Relatively organised approach to data (crystallography data are highly structured)
- Convention is to share derived or reduced data, access to raw data is rare
- Crystallography Information File (CIF) is a de facto exchange standard
  - Maintained by International Union of Crystallography (IUCr)
- Heterogeneity in instrumentation and associated software
- Established system for publishing crystallographic data in UK (Cambridge Crystallographic Data Centre-CCDC)
- Other major databanks
  - Germany (inorganic molecule database)
  - Canada (metals database)
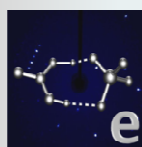  - US (Protein Data Bank -PDB)

# Community & Current Practice (2)

- Publishing datasets
  - Alongside journal articles through publisher mandates
  - Researchers often wish to retain exclusive use of their data
  - Lack of career rewards with respect to data creation and publishing
- Smaller projects at greatest risk
  - Sometimes CIF retained but raw data discarded
  - Data often stored on DVDs or laptops
  - Distributed, local storage -shortage of local curation expertise
  - Quality of metadata for datasets is variable
- Open access
  - eCrystals Federation Project
  - CrystalEye
  - ReciprocalNet (US, Australia, UK)
  - Crystallography Open Database (COD)
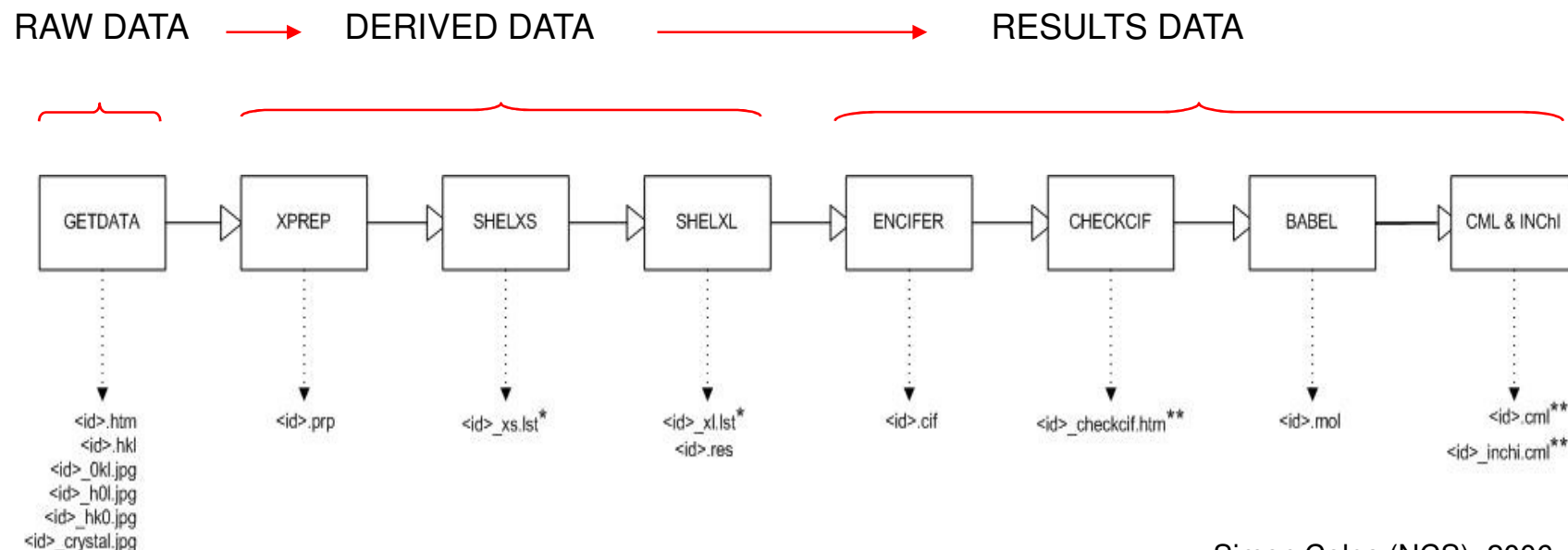  - Chemistry Central (open access publisher)

# Building the eCrystals Repository

- Phenomenal growth in amount of data generated from experiments
  - 40 years ago a PhD student would determine 2-3 structures for a thesis; this can now be easily achieved in a single day
- Only a small proportion is widely and easily accessible
  - Estimated that < 50% of crystal structures are published [Allen 2004]
  - Current data publication process is a bottleneck
- eBank-UK Project
  - JISC funded; three phases Sept. 2003-June 2007
  - UKOLN (lead), University of Southampton, University of Manchester
- eCrystals data repository
  - Open access and rapid dissemination of derived and results data from crystallography experiments
  - Repository platform: ePrints.org software V3
  - Supported by learned society (IUCr) and subject repository (CCDC)
- Linking research data to publications and scholarly communication
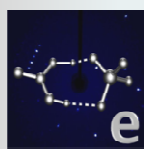- Metadata harvesting and aggregation (OAI-PMH)

**eCrystals** Federation

# EPSRC NCS Crystal Structure Determination Workflow



RAW DATA → DERIVED DATA → RESULTS DATA

GETDATA → XPREP → SHELXS → SHELXL → ENCIFER → CHECKCIF → BABEL → CML & INChI

<id>.htm
<id>.hkl
<id>_0kl.jpg
<id>_h0l.jpg
<id>_hk0.jpg
<id>_crystal.jpg

<id>.prp

<id>_xs.lst*

<id>_xl.lst*
<id>.res

<id>.cif

<id>_checkcif.htm**

<id>.mol

<id>.cml**
<id>_inchi.cml**

Simon Coles (NCS), 2006

- Initialisation: mount new sample
- Collection: collect data
- Processing: process and correct images
- Solution: solve structures

- Refinement: refine structure
- CIF: produce Crystallographic Information File
- Validation: chemical & crystallographic checks
- Report: generate Crystal Structure Report
- CML, INChI

**eCrystals** Federation

# eCrystals Data Repository:
# Example Crystal Structure Report

# Linking Data to Publications

# The Scholarly Knowledge Lifecycle



Both research and learning are cyclical processes
- Research outputs feed into and contribute to knowledge
- Research outputs are based on continuous use and reuse of data i.e. derivative in nature

# Resource Discovery & Reuse

- Simple Dublin Core
  - Crystal structure
  - Title (Systematic IUPAC Name)
  - Authors
  - Affiliation
  - Creation Date
- Qualified Dublin Core (for additional chemical metadata)
  - Empirical formula
  - International Chemical Identifier (InChI)
  - Compound Class and Keywords
- Application Profile: http://www.ukoln.ac.uk/projects/ebank-uk/schemas/
- DOI links: http://dx.doi.org/10.1594/ecrystals.chem.soton.ac.uk/145
- Rights & Citation: http://ecrystals.chem.soton.ac.uk/rights.html

**eCrystals** Federation

# Scaling Up: Towards a Federation

**Interviews, analysis & synthesis:**

> IR Policy & Practice, Laboratory Practice & Workflows, Technical Interoperability & Standards, Metadata Schema & Application Profiles, Semantic Interoperability, Data Citation, Identifiers & Linking, Federation Architectures & Third Party Services, Rights & Licensing, Data Quality & Validation, Preservation, Curation & Sustainability

**Selected Issues (& Recommendations):**

– Diverse laboratory practice
– Instrument manufacturers have proprietary formats
– Data policy needs to reflect laboratory practice
– Data quality criteria and validation (access to raw data)
– Repository must provide control over timing of public visibility-"prior publication" problem
– No disciplinary preservation model

UKOLN

### Scaling Up: Towards a Federation of Crystallography Data Repositories

**Document details**

| Author: | Liz Lyon, Simon Coles, Monica Duke, Traugott Koch |
| --- | --- |
| Date: | 12th May 2008 |
| Version: | 1 0 Final |
| Document Name: | ebank-phase3-report-final.doc |
| Notes: | |

# Data Curation & Preservation

eBank-UK Phase 3: "A Study of Curation and Preservation issues in the eCrystals Data Repository and proposed Federation", Sept. 2007

– Development of preservation strategies and policies

– Audit and certification issues (TRAC, DRAMBORA, NESTOR, ISO International repository audit and certification BOF Group)

– OAIS and Representation Information for crystallography data

– eBank-UK Application Profile and preservation metadata

– e-Prints.org repository platform

A study of Curation and Preservation Issues in the eCrystals Data Repository and Proposed Federation

**eBank-UK Phase 3: WP4**
September 2006 - June 2007

Final Version (Revised): 7th September 2007

Manjula Patel
UKOLN, DCC
University of Bath, UK

Simon Coles
National Crystallography Centre
University of Southampton, UK

**eCrystals** Federation

# Data Curation & Preservation: Recommendations

- Develop a preservation and curation strategy and formal policies to indicate levels of service (e.g. deposit, ingest, validation, dissemination)
- Promote community-supported sustainability plan
- Self-assessment using DRAMBORA toolkit
  - Implement regular audits e.g. annually
  - Produce documentary evidence of compliance
- Maintenance and open access of critical file formats and software
  - Crystallography Information File (CIF)
  - Work-up software e.g. XPREP; SHELX{S,L}; ENCIFER; checkCIF, BABEL
  - Advocate export of raw data from instrumentation as IMG CIF
- Capture relevant Representation Information
- Capture preservation metadata (e.g. versioning; provenance)
  - OAIS Preservation Description Information
  - PREMIS Data Dictionary
  - Extend or augment eBank Metadata Application Profile
- Obtain consensus on Metadata Application Profile
- Seek to automate metadata generation, extraction and maintenance

# Building a Federation of Repositories

# eCrystals Federation Project

- eCrystals Federation Project, Nov 2007 – Mar 2009
- Builds on eBank-UK Phase 3 results
- Led by the UK National Crystallography Service (University of Southampton) with core partners at UKOLN (University of Bath), the Digital Curation Centre and the Unilever Centre (University of Cambridge) – currently 14 supporting partners.
- Integrate and embed open data repository approach into current research practice by engaging data centres, librarians, researchers, publishers and third party information providers
- Harmonise Federation metadata application profile
- Investigate aggregation issues arising from harvesting metadata from Federation repositories
- Enable the Federation of institutional repositories to interoperate with international subject archives (IUCr and CCDC) and other third party harvesters
- Develop approaches to preservation and curation of scientific data in open repositories

# Federation Interoperability

- Roll-out in 2 phases led by University of Southampton
  - Universities Sydney, Drexel, Birmingham, Newcastle with eprints.org platform
  - University Cambridge, STFC, ReciprocalNet, ARCHER with other platforms
  - Establish Federation policies, metadata application profile etc.
- Bi-directional links with derived articles in "publisher repositories", IUCr, RSC, Chemistry Central
- StORe middleware -linking "source" and "output" repositories
- CLADDIER –linking data to publications
- OAI-ORE (Open Archives Initiative – Object Reuse and Exchange)
  - Enable distributed repositories to fully describe and exchange content
  - MicroSoft eChemistry Project

**eCrystals** Federation

# Some challenges

- Data management plans
- Dealing with diverse laboratory practice and workflows
- Appraisal and selection
- Data provenance, audit, tracking
- Citations and versions –persistent identifiers
- Granularity of citations: dataset or values within a dataset
- Instrumentation –proprietary formats
- Access to raw data files for mining and quality control purposes
- Preservation beyond "data" e.g. workflows, blogs, discourse
- Linking across disciplines and sectors
- Collaborative social networks; also "citizen science"
- Semantic integration –controlled vocabularies, ontology etc.

# Selected References

- Liz Lyon, **eBank UK: Building the links between research data, scholarly communication and learning**, ARIADNE, July 2003
- F. H. Allen, **High-throughput crystallography: the challenge of publishing, storing and using the results.** Crystallography Reviews, 10, pp3-15, 2004
- Monica Duke, Michael Day, Rachel Heery, Leslie A. Carr, Simon J. Coles **Enhancing access to research data: the challenge of crystallography** JCDL 2005 Digital Libraries: Cyberinfrastructure for Research and Education, Denver, Colorado, USA June 7-11, 2005
- Scott Carlson, **Lost in a Sea of Science Data,** The Chronicle of Higher Education, June 2006
- Simon Coles, Jeremy Frey and Andrew Milstead **Curation of Chemistry from Laboratory to Publication** UK e-Science All Hands Meeting 2006, Nottingham, UK, 18-21 September 2006
- Liz Lyon, **Dealing with Data: Roles, Rights, Responsibilities and Relationships,** JISC Consultancy Report, June 2007
- Manjula Patel and Simon Coles **A Study of Curation and Preservation issues in the eCrystals Data Repository and proposed Federation,** September 2007
- Liz Lyon, Simon Coles, Monica Duke, Traugott Koch **Scaling Up: Towards a Federation of Crystallography Data Repositories**, May 2008
- **To Share or not to Share, Publication and Quality Assurance of Research Data Outputs,** A Report commissioned by the Research Information Network (RIN), Annex: detailed findings for the eight research areas, June 2008
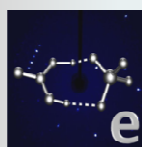
eCrystals Federation

# Thanks …

**…for your attention**

**…to**

Simon Coles:  EPSRC National Crystallography Service,
School of Chemistry, University of Southampton

Liz Lyon: UKOLN, University of Bath

## Questions?

Manjula Patel

UKOLN

University of Bath, UK

http://www.ukoln.ac.uk/

**eCrystals** Federation