# *Institutional Data Repositories for Chemistry*

Simon Coles

School of Chemistry,

University of Southampton, U.K.

s.j.coles@soton.ac.uk

# Why?
# Funding Body Viewpoint

## RESEARCH COUNCILS UK — Together in research

### News release

28 June 2005

Research Data

8. RCUK also notes that one of the benefits of digitisation and publication in digital formats is the ability to provide access to primary research data alongside the traditional article; and it shares the Select Committee's and the Government's view that the data underpinning the published results of publicly-funded research should be made available as widely and rapidly as possible. For a number of years, Research Councils including the AHRB, ESRC and NERC have funded data centres and services which are responsible for preserving, managing and providing access to research data; and these Councils have well-established policies and procedures for preservation and access. CCLRC is currently leading cross-Council consideration of how policy and practice need to be developed with regard to the curation of the data created through the research projects they support. *Further work is needed to develop a common framework of policies and procedures for determining what sets of data are collected, whether in university or in Council-run repositories or elsewhere; and how and on what terms they are made accessible to the research community and others*

# Why?
# Curation in the Laboratory

"Data from experiments conducted as recently as six months ago might be suddenly deemed important, but those researchers may never find those numbers – or if they did might not know what those numbers meant"

"Lost in some research assistant's computer, the data are often irretrievable or an undecipherable string of digits"

"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."
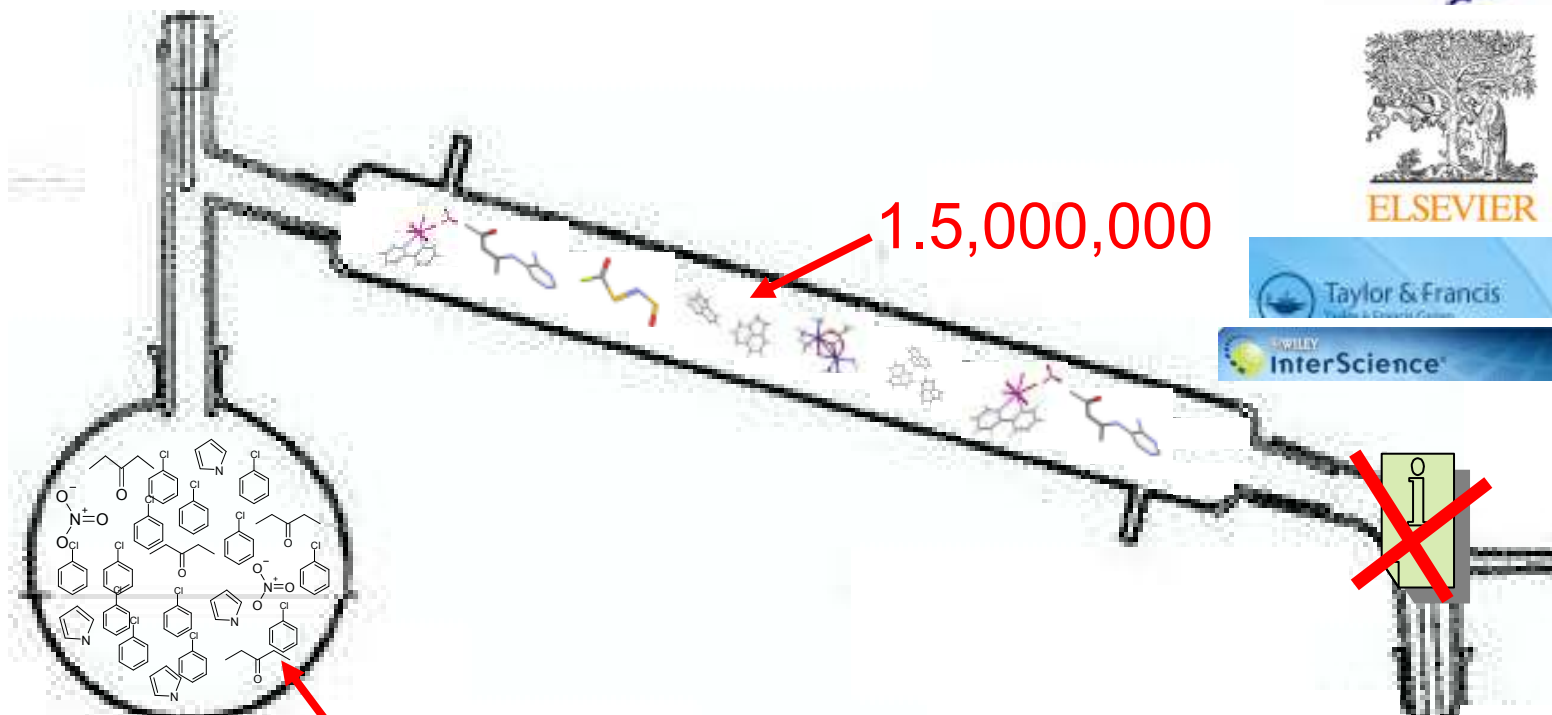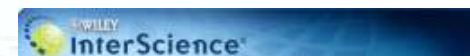
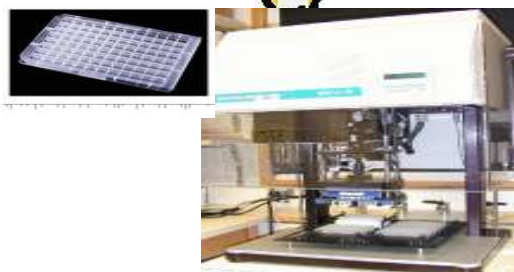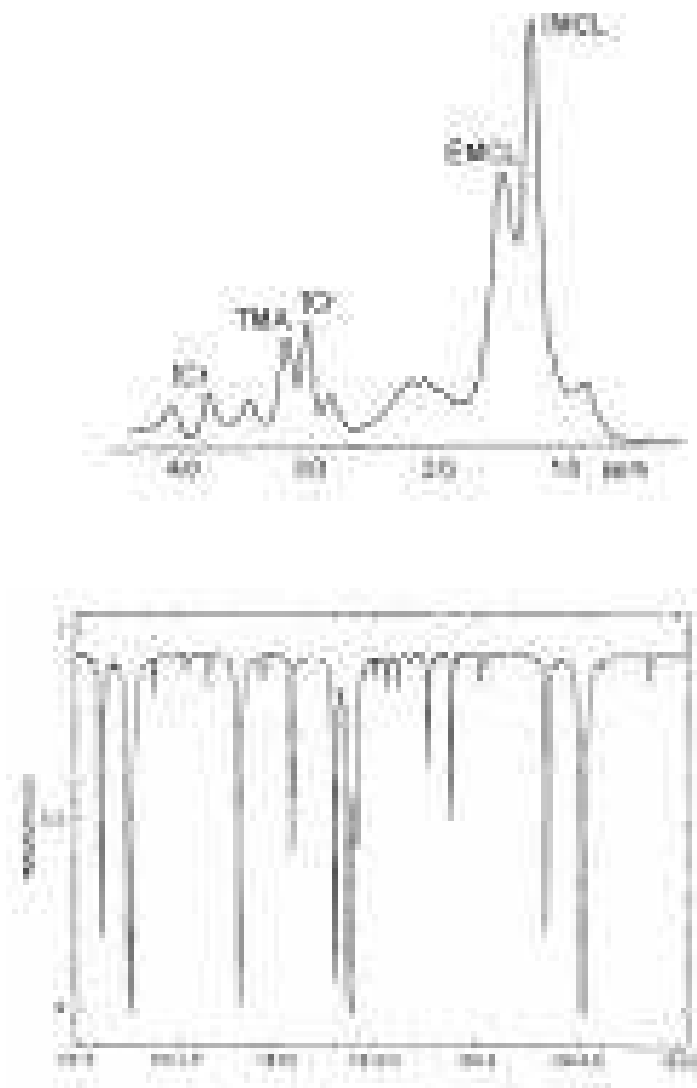'Lost in a Sea of Science Data' S.Carlson, The Chronicle of Higher Education (23/06/2006)

1,5,000,000

30,000,000

450,000

# Why?
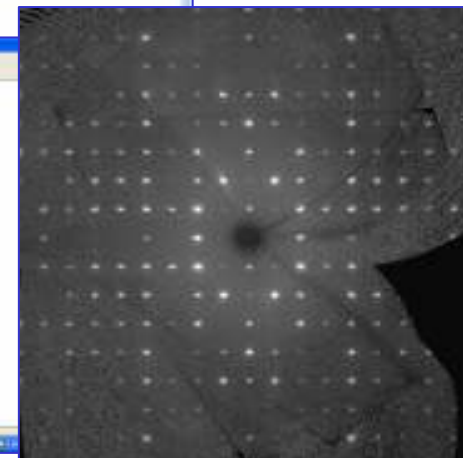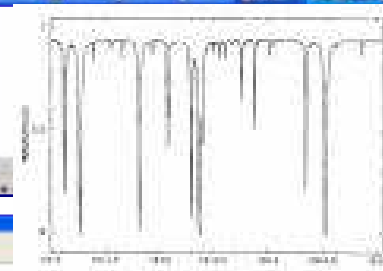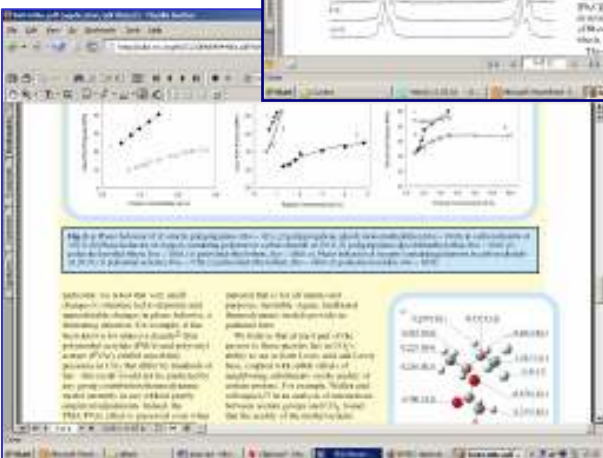## Publishing Data and Information Loss

# Separating Data from Interpretations

Intellect & Interpretation (Journal article, report, etc)

Underlying data (Institutional data repository)

# Data capture and curation at the point of generation in the laboratory

# The Repository for the Laboratory – R4L

# Laboratory IRs and Information Management



LABORATORY EQUIPMENT

LABORATORY REPOSITORY

#1 INGEST PROCESS

EXTERNAL REPOSITORIES

OAI

PRIORITY ASSERTION SERVICE

DATA REUSE

#1 Collaborate with instrument manufacturers to develop protocols for data deposition & metadata

#2 Devise a service to establish a reliable timestamp to provide a legally sound guarantee of priority

#3 Develop management protocols and tools to manage heterogeneous and multiple datasets in a repository

#4 Develop a tool to generate a formal description of the experimental process and compare data from different analyses

#5 Collaborate with ALPSP and the eBank-UK project to develop data citation and aggregation protocols

DATA DISSEMINATION AND AGGREGATION (EBANK-UK PROJECT)

#3 REPOSITORY MANAGEMENT

OAI

#4 SCIENTIFIC DATA REPORT

OAI-PMH

INSTITUTIONAL REPOSITORY

OAI

#5 DATA CITATION REPORTING

ARTICLE

# The R4L Repository

Create new compound

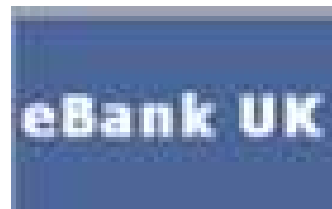Add experiment data and metadata

Deposit

Search / Browse

# Data dissemination and curation by the scientist and host institution

# eBank-UK and the eCrystals Repository

# The eCrystals Data Archive

**University of Southampton — Crystal Structure Report Archive**

- Home
- About
- Browse
- Search
- Register
- User Area
- Help

## 6,7,9,10,12,13,15,16-Octahydro-benzo-1,4,7,10,13-pentaoxacyclopentadecin

Simon J Coles, Michael B Hursthouse,
Jeremy G Frey and Esther Rousay.

University of Southampton

$C_{14}H_{20}O_5$

InChl=1/C14H20O5/c1-2-4-14-13(3-1)18-11-9-16-7-5-15-6-8-17-10-12-19-14/h1-4H,5-12H2

**DOI:** 10.594/ecrystals.chem.soton.ac.uk/145
**Compound Class:** Organic
**Keywords:** crown ethers
**Creation Date:** 07 October 2004
**Deposited By:** A.N. Admin
**Deposited On:** 20 February 2006

**Available Files**

### Depositor Comments

Structure already known, but accurately redetermined for a local research project.

### Data collection parameters

| Chemical formula | C14 H20 O5 |
|---|---|
| Crystallisation Solvent | |
| Crystal morphology | Plate |
| Crystal system | Orthorhombic |
| Space group symbol | Pbca |
| Cell length a | 16.4963(18) |
| Cell length b | 8.325(3) |
| Cell length c | 20.061(6) |
| Cell angle alpha | 90.00 |
| Cell angle beta | 90.00 |
| Cell angle gamma | 90.00 |
| Data collection temperature | 120(2) |

### Refinement results

| Solution figure of merit | 0.0409 |
|---|---|
| R Factor (Obs) | 0.0487 |
| R Factor (All) | 0.0977 |
| Weighted R Factor (Obs) | 0.1008 |
| Weighted R Factor (All) | 0.1192 |

Citations: Coles, S.J., Hursthouse, M.B., Frey, J.G. and Rousay, E. (2004), Southampton, UK, University of Southampton, Crystal Structure Report Archive. (doi:10.1594/ecrystals.chem.soton.ac.uk/145)

### Final Result

| 04sjc0831.cif | 13k |
|---|---|
| 04sjc0831.cml | 6k |

### Validation

| 04sjc0831_checkcif.htm | 7k |
|---|---|

### Refinement

| 04sjc0831.res | 6k |
|---|---|
| 04sjc0831_pl.lst | 34k |

### Solution

| 04sjc0831.prp | 6k |
|---|---|
| 04sjc0831_xs.lst | 39k |

### Processing

| 04sjc0831.hkl | 702k |
|---|---|
| 04sjc0831.htm | 10k |
| 04sjc0831_0kl.jpg | 57k |
| 04sjc0831_h0l.jpg | 85k |
| 04sjc0831_hk0.jpg | 88k |

### Data Collection

| 04sjc0831_crystal.jpg | 17k |
|---|---|

### Other Files

| 04sjc0831.doc | 78k |
|---|---|
| 04sjc0831.fcf.txt | 155k |

http://ecrystals.chem.soton.ac.uk

# Metadata Publication

- Using simple Dublin Core
  - Crystal structure
  - Title (Systematic IUPAC Name)
  - Authors
  - Affiliation
  - Creation Date
- Additional chemical information through Qualified Dublin Core
  - Empirical formula
  - International Chemical Identifier (InChI)
  - Compound Class & Keywords
- Specifies which 'datasets' are present in an entry

- DOI *http://dx.doi.org/10.1594/ecrystals.chem.soton.ac.uk/145*

- Rights & Citation *http://ecrystals.chem.soton.ac.uk/rights.html*

- Application Profile *http://www.ukoln.ac.uk/projects/ebank-uk/schemas/*

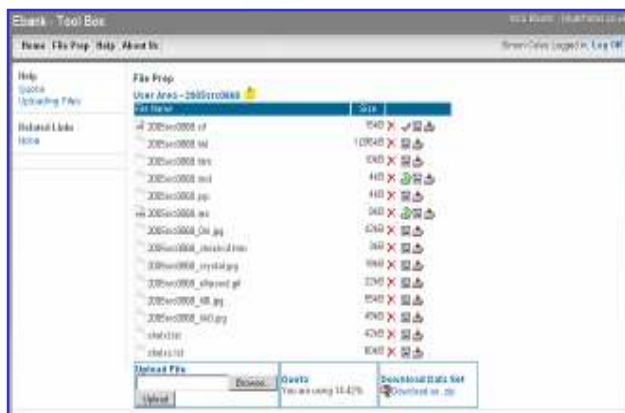# Metadata and Data Quality Control

## Data manipulation toolbox

## Associated Metadata

Value added

Format conversion

# Laboratory Data Management and Archive



University of Southampton
**Crystal Structure Report Archive**

Home
About
Browse
Search
Register
User Area
Help

## 6,7,9,10,12,13,15,16-Octahydro-benzo-1,4,7,10,13-pentaoxacyclopentadecin

**Origination:** Esther Rousay and Jeremy G Frey.

**Data Collection:** Simon J Coles.

**Structure Determination:** Simon J Coles and Micheal B Hursthouse.

University of Southampton

$C_{14}H_{20}O_5$

InChl=1/C14H20O5/c1-2-4-14-13(3-1)18-11-9-16-7-5-15-6-8-17-10-12-19-14/h1-4H,5-12H2

**Compound Class:** Organic
**Keywords:** Benzo-15-crown-5
**Creation Date:** 07 October 2004

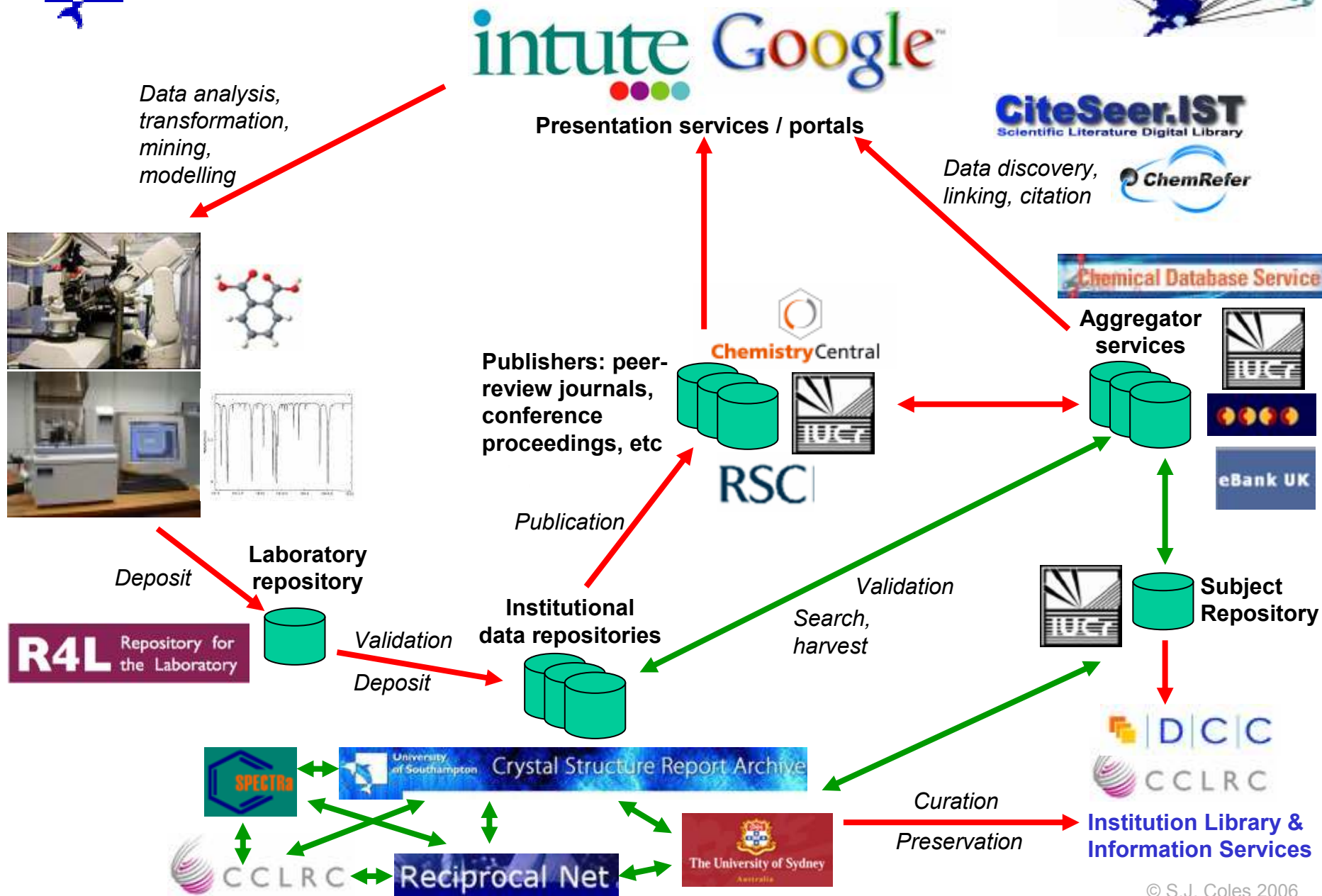**Available Files**

**Final Result**

# Institutional data repositories and harvesting, aggregation and curation by data centres and third party services

## eBank-UK Phase 3 – The eCrystals Federation

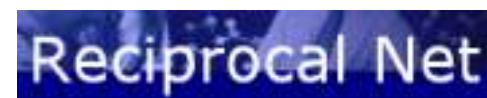# The eCrystals 'Global Federation' Model

intute Google

**Presentation services / portals**

CiteSeer.IST
Scientific Literature Digital Library

ChemRefer

*Data analysis, transformation, mining, modelling*

*Data discovery, linking, citation*

Chemical Database Service

ChemistryCentral

**Aggregator services**

IUCr

**Publishers: peer-review journals, conference proceedings, etc**

RSC

eBank UK

*Publication*

*Validation*

*Search, harvest*

**Laboratory repository**

*Deposit*

R4L Repository for the Laboratory

*Validation*

*Deposit*

**Institutional data repositories**

IUCr

**Subject Repository**

DCC
CCLRC

University of Southampton Crystal Structure Report Archive

SPECTRa

*Curation*

*Preservation*

**Institution Library & Information Services**

CCLRC ↔ Reciprocal Net ↔ The University of Sydney Australia

# Exploring the heterogeneous landscape of data repositories

Crystal Structure Report Archive

- Different software platforms

- Different administrative domains

- Different Institutional structure

- Institutional vs Subject repositories
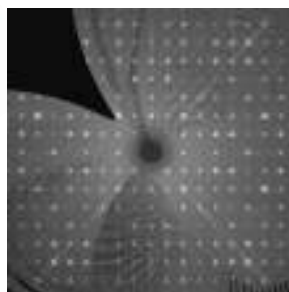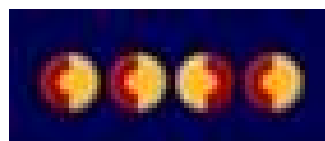
- Data Repository Interoperability    ORE

# Preservation and curation
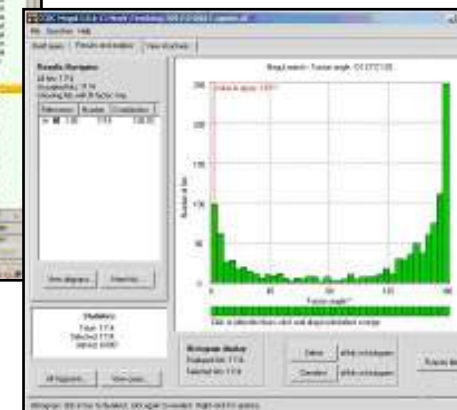# by data centres & Institutions



CCLRC

G bytes

M bytes

CASPAR  D|C|C

k bytes

IUCr

**Institution Library & Information Services**

# Harvesting, aggregation, value addition and curation by data centres



Cambridge Crystallographic Data Centre

Chemical Database Service

eBank UK

# The relationship with (conventional?) publication protocols and procedures

- Discipline-based publication

- Domain-based publication

- Open Access publication

# Aggregation, linking and information provision by third party services

• Indexing and aggregating with other datasets

• Aggregating and linking between datasets and articles

• Integration into information portals