

Jim Downing, University of Cambridge

Dr. CrystalEye

Or:
How
I Learned
To
Stop
Worrying
And
Love
The
Web

From Desktop to Data Repository



I'm going to tell you a story about a PhD in my research group, and a data system he developed, called CrystalEye. Like all good stories, there's a meta-story between the lines. Since the story teller gets to pick the meta-story, today it's about repositories.



In the beginning...

Nick Day started his PhD under the supervision of Peter Murray-Rust. Quantum computational programs to calculate molecular structures. Compare calculated structures to actual structures, as determined by X-Ray crystallography to work out when and why the programs got it wrong. As things transpired, the collection and publication of the X-Ray crystallography turned out to an interesting area in itself – and it's this side of Nick's PhD that interests us most today.

How to get Open crystallographic structure data (2005)



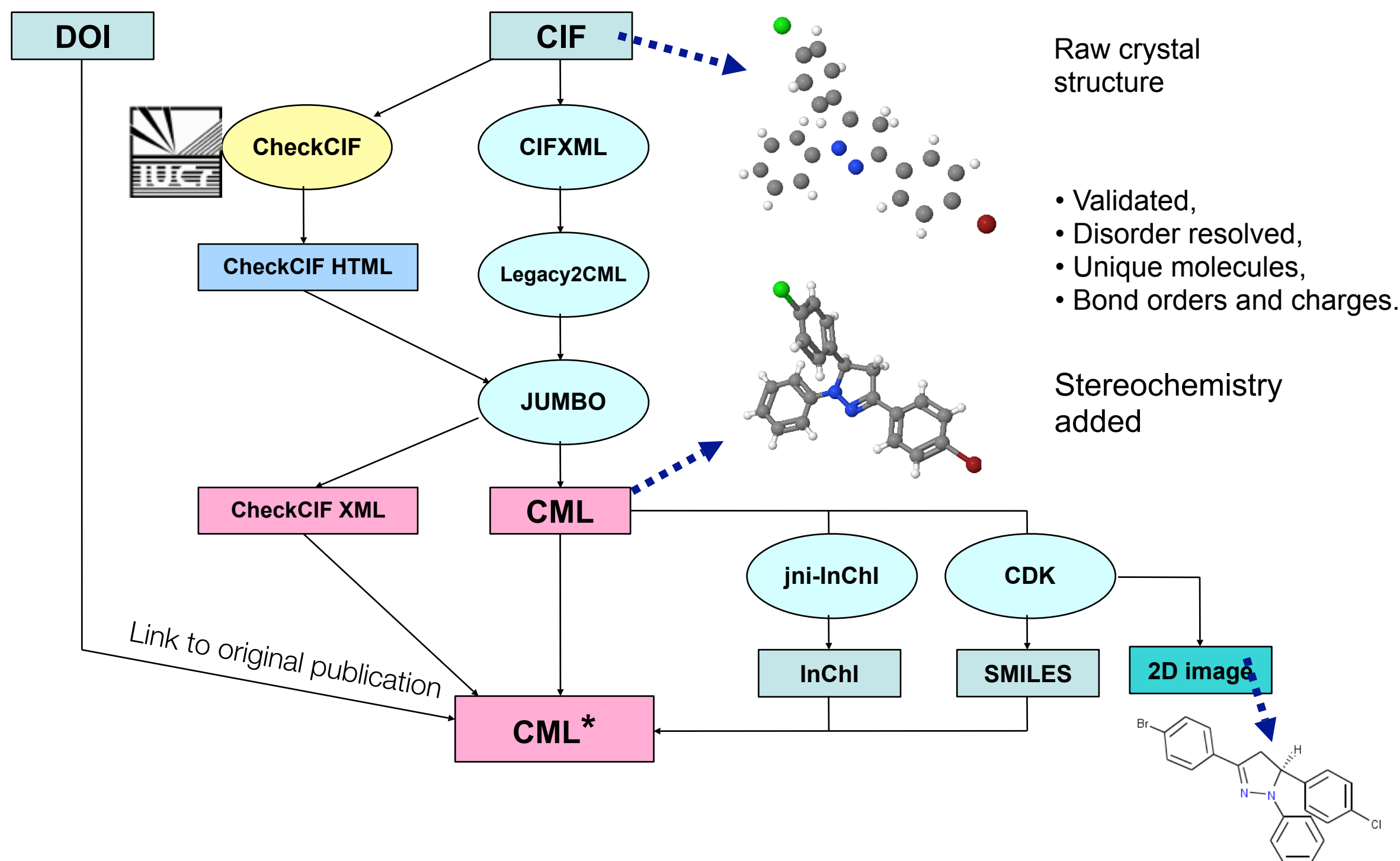
Nick needed large amounts of open structure information to compare the computational outputs with.

There was no database of open crystallographic data.

There is crystallographic data out there – on journal websites. The information comes from the websites of acta journals (especially the Acta Crystallographica family published by the IUCr) that specialize in reporting X-Ray structure determinations, and from the supporting information of other chemistry journal publications.

Nick wrote a web spider to find it and collect it.

CrystalEye Data Processing

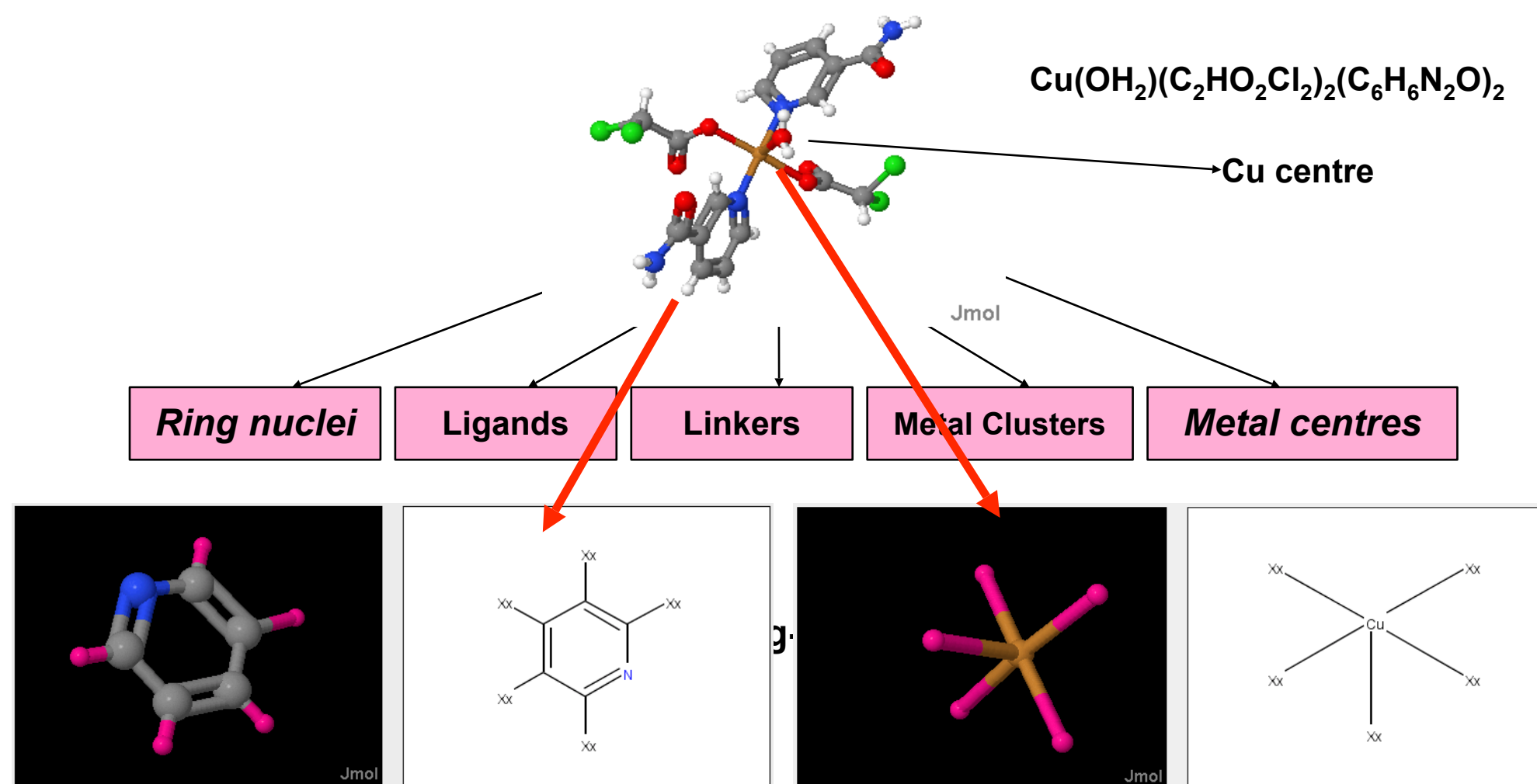


Nick needed to convert his data into Chemical Markup Language, an XML for chemical data. Being a good geek, he didn't miss the opportunity follow some interesting side lines, including developing heuristic approaches to fix and enhance missing data.

(This slide shows the processing of the crystallographic data once it has been collected by the spider. The detail isn't necessary to follow the story)

CrystalEye Data Processing 2

- Automatic generation of fragments



... and break apart molecules to form a fragment library (these can be used for an empirical approach to predicting 3D molecule structures). Nick also had web authoring skills, so he created HTML pages for the data he was collecting as part of the processing. To help himself keep tabs on the growing collection, he created a small but growing website on his desktop machine.



UNIVERSITY OF
CAMBRIDGE

Crystal

Home

Search

Browse Issues

RSS feeds

Bond Lengths

Greasemonkey

FAQ

Acta Crystallographica - Section E:

- 2008
 - [Issue 03-00](#)
 - [Issue 02-00](#)
 - [Issue 01-00](#)
- 2007
 - [Issue 12-00](#)
 - [Issue 11-00](#)
 - [Issue 10-00](#)
 - [Issue 09-00](#)
 - [Issue 08-00](#)
 - [Issue 07-00](#)
 - [Issue 06-00](#)
 - [Issue 05-00](#)
 - [Issue 04-00](#)
 - [Issue 03-00](#)

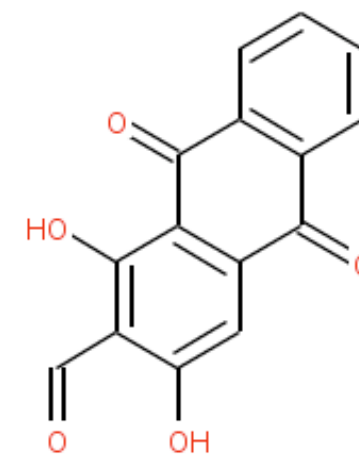
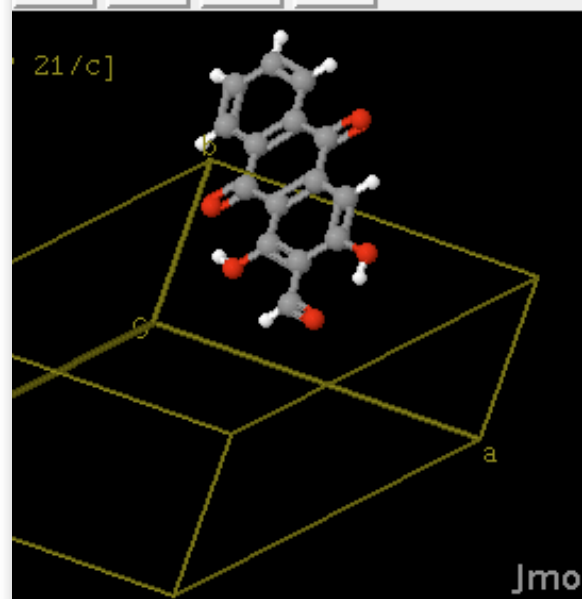
ACTA CRYSTALLOGRAPHICA
SECTION E, 2008, ISSUE 03-00

ORGANIC STRUCTURES

Published Formula (clickable)	Article	Summary
C₁₅H₈O₅	view	view
C₁₄H₂₁N₃O₂S	view	view
C₁₉H₂₁NOS	view	view
C₁₆H₁₅F₂NO₃	view	view
C₁₈H₁₆N₂O₃	view	view
C₁₇H₂₈O₅	view	view
C₉H₁₁FN₂O₃	view	view
C₉H₈N₂O₂	view	view
C₁₂H₁₃NO₅	view	view
C₁₆H₁₈N₄O₇S	view	view
C₁₉H₁₇ClN₄O₂S	view	view

<<- <- -> ->>

Prev 1/1 Next



Browse

Drill Down By Journal, By Issue,
and then by structural sub-
components

1,3-Dihydroxy-9,10-dioxo-9,10-dihydroanthracene-2-carbaldehyde

[OPEN](#) [DATA](#)

[<< Table of Contents](#)

Publisher: Acta Crystallographica

Journal: Section E

Year/Issue: 2008/03-00

Article (via DOI): [10.1107/S1600536808004169](https://doi.org/10.1107/S1600536808004169)

Compound Class: organic

Date Recorded: 2008-01-28

Contact Author: Retailleau, Pascal

e-mail: pascal.retailleau@icsn.cnrs-gif.fr

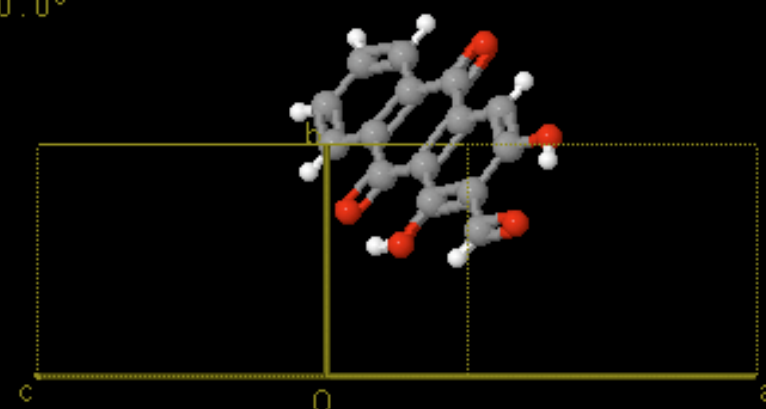
Data collection parameters

Chemical formula sum	C ₁₅ H ₈ O ₅
Chemical formula moiety	C ₁₅ H ₈ O ₅
Crystal system	monoclinic
Space group H-M	P 21/c
Space group Hall	-P 2ybc
Data collection temperature	293.0

Refinement results

R Factor (Obs)	0.051
R Factor (All)	0.081
Weighted R Factor (Obs)	0.132
Weighted R Factor (All)	0.15

```
-P 2ybc [P 21/c]  
a=10.547Å  
b=5.669Å  
c=20.231Å  
α=90.0°  
β=110.6°  
γ=90.0°
```



Jmol

Show no. of unit cells along axis:

a:
b:
c:

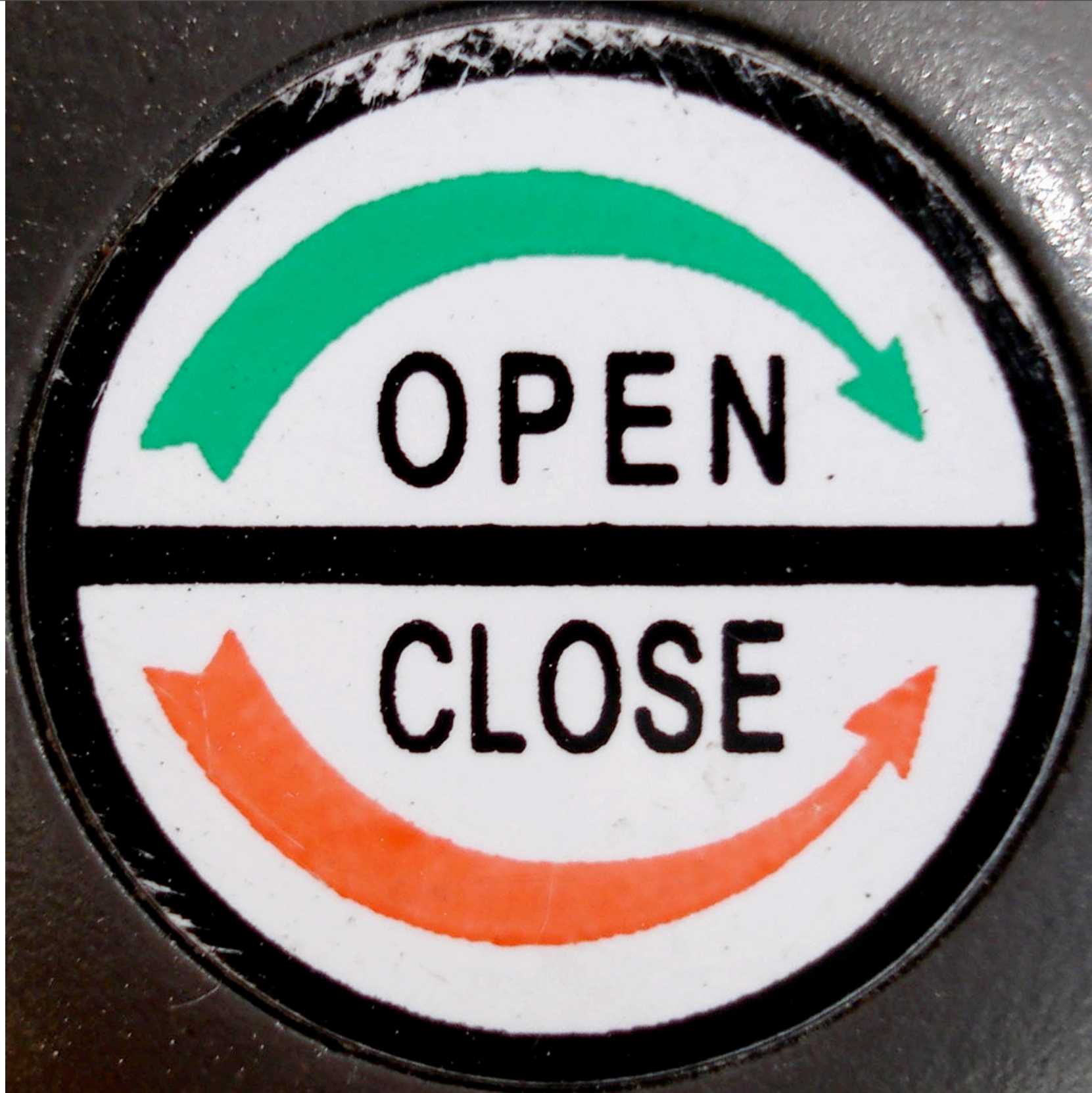
Enter Jmol script:

View Record

Metadata specific to data domain

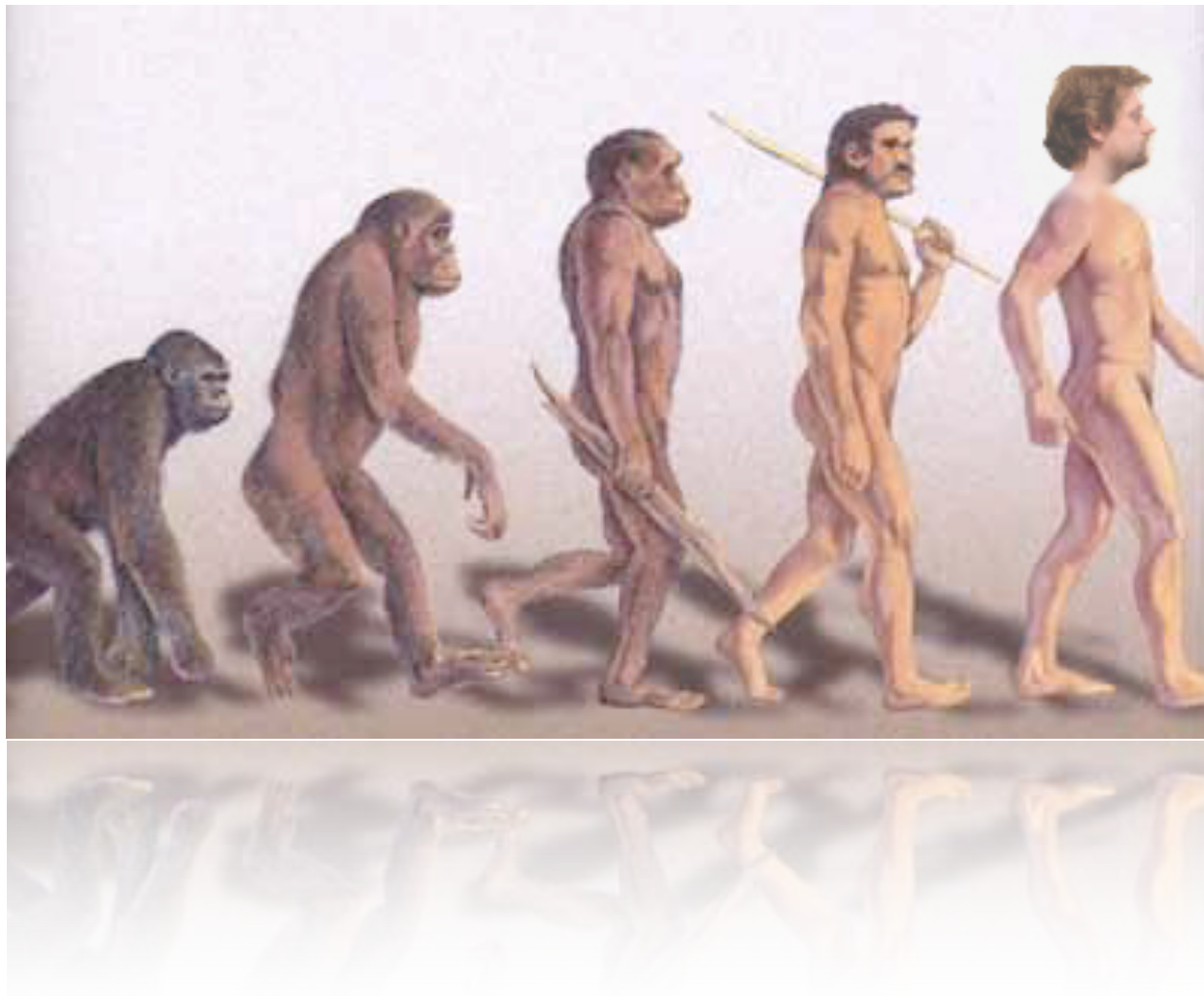
The applet controls or the script loading input can be used to rotate, zoom and otherwise alter and manipulate the visualisation.

.. and so Nick had built himself a web site of his data on his desktop.



Open Data

So Nick had a growing pile of result data, and a visualisation system of Java classes that generated CML and HTML web pages. Peter Murray-Rust and Nick decided to share his results with the community. For myself as part-time system administrator this was no problem – deployment just meant a couple of cron jobs, a large file system partition and line or two of apache config.



Evolution

Sharing your data hopefully means getting users. And having users usually means getting good ideas about what to develop next. As the data set grew, Nick evolved CrystalEye, adding features to make this data collection more useful.

CrystalEye (beta)

Search:

1: Search for substructures in CrystalEye using a SMILES string.

Enter SMILES:

c1ccccc1

2: Search for crystals by cell parameters

Enter cell parameters:

a ± Å
b ± Å
c ± Å
α ± °
β ± °
γ ± °

Search

CrystalEye (beta)

Searching...

Search results for SMILES string "c1ccccc1" (in reverse chronological order) :

Showing results 1 to 25 of 1793

1: c1ccc(cc1)P(c1ccccc1)CCB1(C2CCCC1CCC2)=P(c1ccccc1)(c1ccccc1)CCB1(C2CCCC1C(c1ccccc1)CCBC12C3(CCC1)C2CC3

[The American Chemical Society, Organometallics, 2008, 5, article om7010886, datablock I](#)

2: C=C1C2(C3CC4C(C(=CC3C14)C(C)C)=O)C1COC(C1CC2)=O.C=C1C2(C3CC4C(C(=CC(C)C)=O)C1COC(C1CC2)=O

[Royal Society of Chemistry, Chemical Communications, 2008, 10, article b718754h, datablock I](#)

3: OC[C@H](CC(C(=O)C)=C1[C@@H]([C@H]2([C@@]3(CC[C@]4[C@@]5(CC[C@H](C[C@@]([C@@H]3(C[C@@H]2(O1))))O)C)C)([CH2]))C

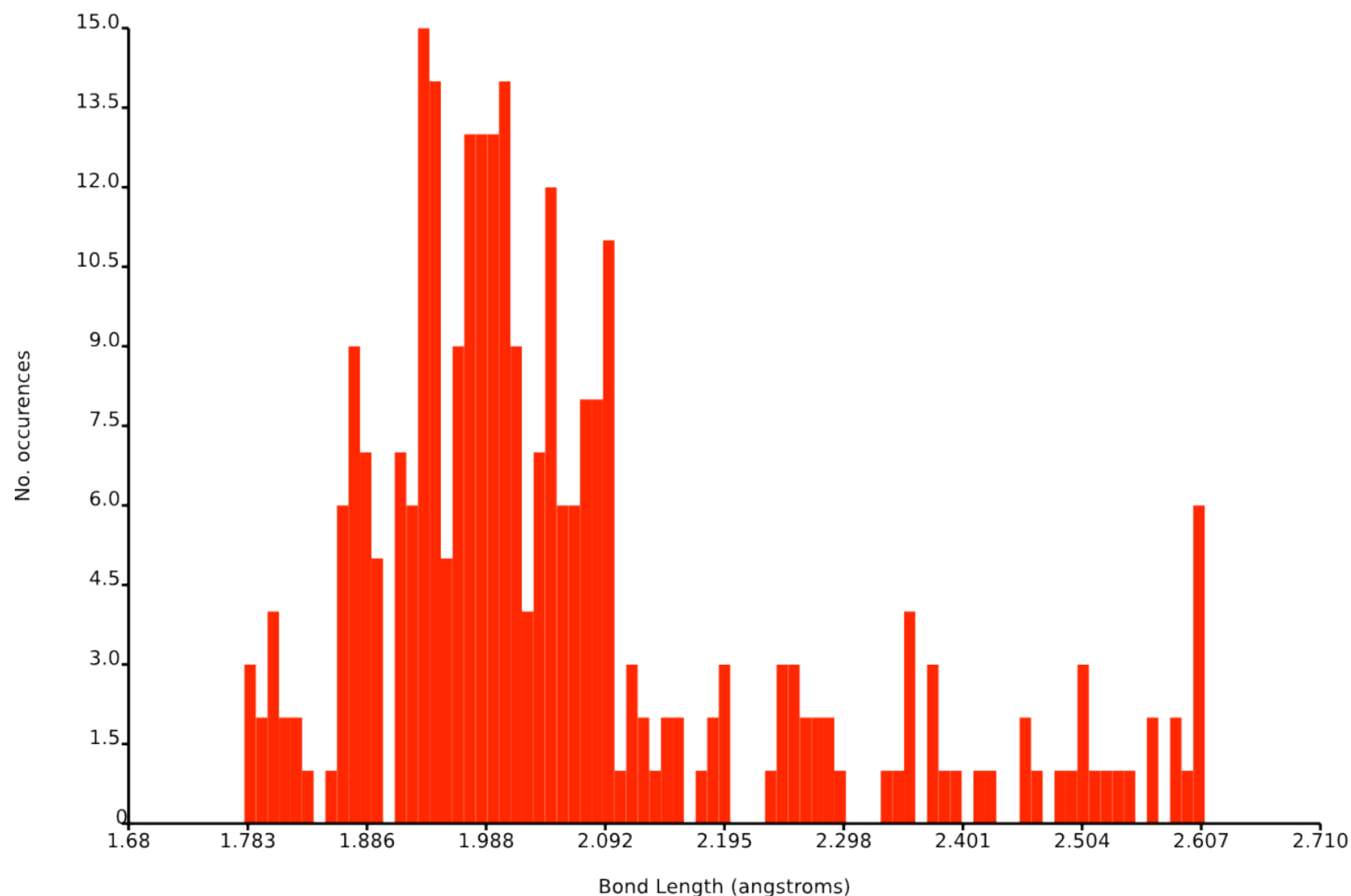
[Acta Crystallographica, Section E, 2008, 03-00, article rz2197, datablock I](#)

Search

The search tool was the first dynamic part of the application, and runs as a simple set of Java servlets.

N.B. that the search is chemistry specific – the first search method might be used for any collection of molecule data, the second is specific to crystallographic data.

C-Cu Bond Lengths in CrystalEye (Last updated 2008-03-20)
for non-disordered, unconstrained atoms in structures where temperature ≤ 200.0 and r-factor ≤ 0.05




Bond Length Histograms

Bond lengths histograms are a good example of the kind of processing and data checking that can only be done with a large collection of data.

Here we see the distribution of bond lengths between copper and carbon atoms in structures in CrystalEye. Nick generates these histograms for every type of bond.

Outliers: bad experimentation, also helped spot edge cases where computational codes failed to take account of some effect.

[Home](#)
[All items \(64\)](#)
[Starred items](#) ★
[Trends](#) 
[Your shared items](#) 
☐ [Friends' shared items](#)
 [Jean-Claude Bradley](#)
[Manage friends »](#)

[Add subscription](#) [Discover »](#)

Show: **updated** - [all](#) [Refresh](#)

☐ [a \(2\)](#)
 [Jimmy Nilsson's weblog \(2\)](#)

☐ [c \(22\)](#)
 [BioMed Central \(1\)](#)
 [DSpace Wiki - Recent ... \(2\)](#)
 [Jon's Radio \(1\)](#)
 [Lifehacker \(16\)](#)
 [O'Reilly Radar \(2\)](#)

☐ [d \(40\)](#)
 [CrystalEye: All Struc...](#)
 [Delicious DSpace \(2\)](#)
 [James Paice's Recent ... \(1\)](#)
 [JL Java Announcements \(8\)](#)
 [Liberal Education Today \(2\)](#)
 [MacOSXHints.com \(5\)](#)
 [Science Blog - \(5\)](#)
 [TechCrunch \(17\)](#)

[Manage subscriptions »](#)

CrystalEye: All Structures

Show: [0 new items](#) - [all items](#)

[Mark all as read](#)

[Refresh](#)

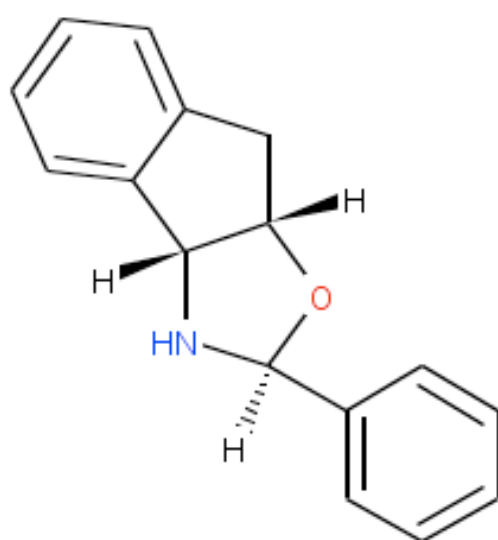
[show details](#)

[Original enclosure](#)

☐ Add star
 ☐ Share
 ☐ Email
 ☒ Mark as read
 ☐ Edit tags: d

★ **Summary page for crystal structure from DataBlock shelxl in CIF b719624esup1 from article b719624e in issue 2008/7 of Royal Society of Chemistry, Organic and Biomolecular Chemistry.** »

Mar 26, 2008 (yesterday)



[Original enclosure](#)

☐ Add star
 ☐ Share
 ☐ Email
 ☒ Mark as read
 ☐ Edit tags: d

★ **Summary page for crystal structure from DataBlock aba2 in CIF b740000esup1 from article b740000e in issue 2008/7 of Royal Society of Chemistry, Organic and Biomolecular Chemistry.** »

Mar 26, 2008 (yesterday)

[Previous item](#)

[Next item](#)

more than 60 items

Feeds for all

Development of Atom feeds allowed embedded HTML and CML as enclosures

There were some problems with the existing CMLRSS feeds – including the CML inline was a big performance hit, and couldn't be done with standard XML libraries because the CML was too large. We started to use Atom, and link to the CML data using enclosures. Using Atom also allowed us to link to images in the entries, which meant the feeds looked good in common-or-garden feed readers.



Harvesting the Atom

RFC 5005
Feed paging
and archiving

OAI-PMH was on the development road map, but users wanted to harvest the data using RSS feeds. RSS feeds aren't effective for CrystalEye harvesting, the data tends to arrive in big clump updates rather than a nice steady trickle, and there was no standard way for harvesters to discover or recover when they missed items. RFC 5005, which was published in its final form last September, extends Atom to fix these problems through a few special elements, and requirements on how the server must maintain and publish archive feed documents. We implemented RFC5005 in our main Atom feed, and published a simple harvester client application.

RDF metadata

1,3-Dihydroxy-9,10-dioxo-9,10-dihydroanthrac

OPEN

DATA

[<< Table of Contents](#)

Publisher: Acta Crystallographica

Journal: Section E

Year/Issue: 2008/03-00

Article (via DOI): [10.1107/S1600502808001600](#)

Compound Class: organic

Date Recorded: 2008-01-28

Contact Author: Retailleau, Pascal

e-mail: pascal.retailleau@icsn.cnrs.fr

Data collection parameters

Chemical formula sum	C ₁₄ H ₈ O ₄
Chemical formula moiety	C ₁₄ H ₈ O ₄
Crystal system	monoclinic
Space group H-M	P2 ₁ /c
Space group Hall	-C2H

```
<rdf:RDF
  xmlns:j.0="http://wwmm.ch.cam.ac.uk/crystaleye/dictionary#"
  xmlns:j.1="http://purl.org/dc/terms/"
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
>
  <rdf:Description rdf:about="http://wwmm.ch.cam.ac.uk/crystaleye/s
    <j.1:contributor>"H\ 'edoux, Alain"</j.1:contributor>
    <j.1:contributor>"Hernandez, Olivier"</j.1:contributor>
    <j.1:contributor>"Masson, Olivier"</j.1:contributor>
    <j.1:contributor>"Lefebvre, Jacques"</j.1:contributor>
    <j.1:contributor>"Guinet, Yannick"</j.1:contributor>
    <j.1:contributor>"Descamps, Marc"</j.1:contributor>
    <j.1:contributor>"Papoular, Robert"</j.1:contributor>
    <j.0:oscarAnnotation rdf:resource="http://wwmm.ch.cam.ac.uk/cry
    <j.0:oscarAnnotation rdf:resource="http://wwmm.ch.cam.ac.uk/cry
    <j.0:oscarAnnotation rdf:resource="http://wwmm.ch.cam.ac.uk/cry
    <j.0:oscarAnnotation rdf:resource="http://wwmm.ch.cam.ac.uk/cry
    <j.0:oscarAnnotation rdf:resource="http://wwmm.ch.cam.ac.uk/cry
    <j.0:_publ_section_title rdf:datatype="http://example.com/json"
    <j.0:_journal_name_full rdf:datatype="http://example.com/json">
    <j.0:location rdf:resource="http://wwmm.ch.cam.ac.uk/crystaleye
    <j.0:location rdf:resource="http://wwmm.ch.cam.ac.uk/crystaleye
    <j.0:location rdf:resource="http://wwmm.ch.cam.ac.uk/crystaleye
    <j.0:AcceptanceDate rdf:datatype="http://www.w3.org/2001/XMLSchema
  </rdf:Description>
  <rdf:Description rdf:about="http://wwmm.ch.cam.ac.uk/crystaleye/1
    <j.0:latitude rdf:datatype="http://www.w3.org/2001/XMLSchema#fl
    <j.0:address>" Laboratoire L\ 'eon Brillouin (UMR 12 CNRS-CEA),
    <j.0:longitude rdf:datatype="http://www.w3.org/2001/XMLSchema#f
  </rdf:Description>
  <rdf:Description rdf:about="http://wwmm.ch.cam.ac.uk/crystaleye/1
    <j.0:address>"LCSIM UMR 6511 CNRS - Universit\ 'e de Rennes 1,
    <j.0:latitude rdf:datatype="http://www.w3.org/2001/XMLSchema#fl
```

Using a link element in the head of the HTML pages, machine clients can obtain machine readable metadata about the structure in RDF format.

CrystalEye query for: Redfern S A T

Get this page as: [Google Earth/Maps format](#)

[An in situ neutron diffraction study of cation disordering in synthetic qandilite Mg₂TiO₄ at high temperature Sample: First series, T = 1020 C \(raw CML\)](#)

published in *American Mineralogist*

Authors:

- [Short S M](#)
- [O'Neill H St C](#)
- [Redfern S A T](#)
- [Kesson S](#)

Chemical named entities

- [Mg₂TiO₄](#)

Chemical ontology terms

- [cation](#)
- [neutron](#)
- [An](#)

[An in situ neutron diffraction study of cation disordering in synthetic qandilite Mg₂TiO₄ at high temperature Sample: First series, T = 1054 C \(raw CML\)](#)

published in *American Mineralogist*

Authors:

- [Short S M](#)
- [O'Neill H St C](#)
- [Redfern S A T](#)
- [Kesson S](#)

Chemical named entities

- [Mg₂TiO₄](#)

Chemical ontology terms

- [cation](#)

Mashing up RDF data

CrystalEye RDF data, concepts extracted, reindexed reindexed by by author, chemical name etc.

Andrew Walkingshaw took the CrystalEye RDF data, added some data of his own using the OSCAR3 chemistry natural language processing application, and mashed it back together to provide views and indices into the CrystalEye data we hadn't even thought of creating. Outside In-novation! It gets much, much sexier than this screenshot, but I can't show you too much as Andrew is presenting his work at XTech next month!

The global distribution of crystallography: papers in CrystalEye, 2000-2007

Dr Andrew Walkingshaw

Unilever Centre for Molecular Science Informatics

for Open Repositories 2008

28th March, 2008

<http://www.lexical.org.uk/>

<http://wwwmm.ch.cam.ac.uk/blogs/walkingshaw/>

Summary of Development

- Crystallographic data is collected daily from publisher websites, processed and enhanced.
- Web resources are constructed as part of data processing and published as static files by Apache httpd
 - Browse & View
 - Feeds for viewing and for harvesting
- Linked Machine Readable Metadata
- Data resources indexed for Search (implemented as Java servlets)
- Aggregate processing and results publication

The Future

Development on and around
CrystalEye continues...



Planned work

- Refactoring!
 - Enhanced, archived atom feeds replacing CMLRSS throughout
 - Decouple the spider, processing and publishing parts
 - More auto-discovery
- Departmental crystallography repository, using CrystalEye and the JISC SPECTRa tools.
- New features
 - Dumpfile download (through S3)
 - SWORD support
 - ORE support (Microsoft ORE project)
 - OAI-PMH support (for e-Crystals federation)



Nearly there!

Don't worry, we're in the final furlong, and I'm about to start talking about repositories...

That's the end of my narrative of the CrystalEye story so far, so I want to move on to talk about the meta-story, which is how CrystalEye is relevant to repositories...

Is CrystalEye a Repository?

It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

Is CrystalEye a Repository?

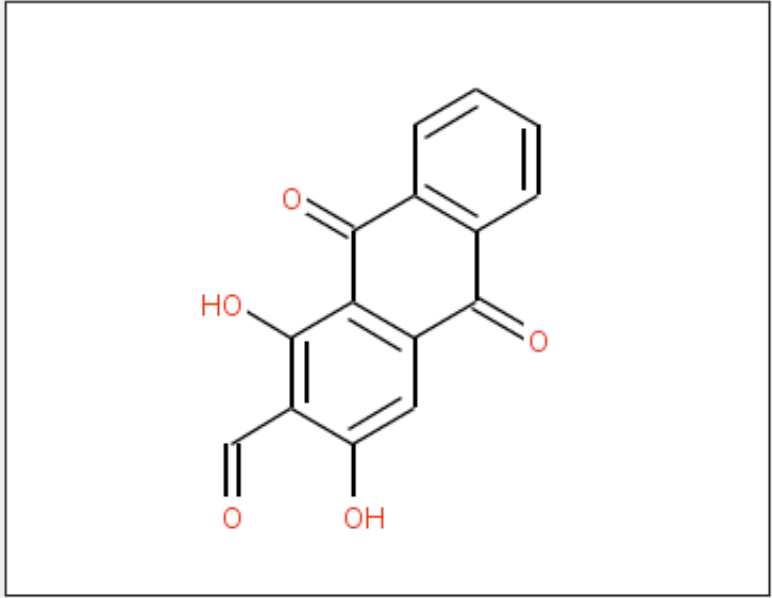
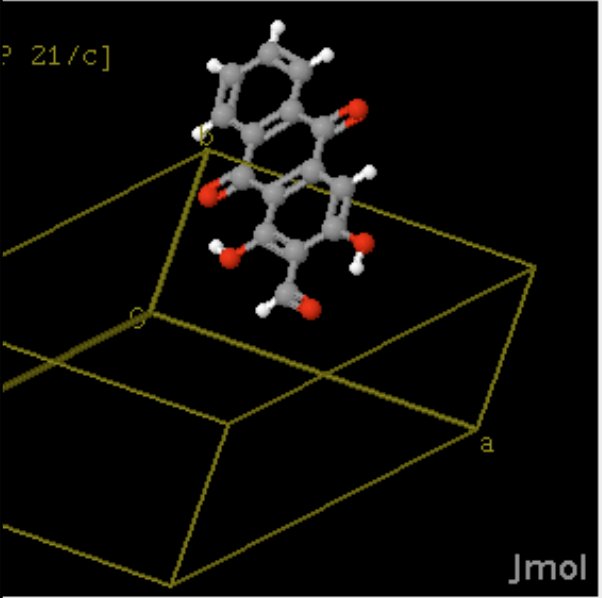
Browse

**ACTA CRYSTALLOGRAPHICA
SECTION E, 2008, ISSUE 03-00**

ORGANIC STRUCTURES

Published Formula (clickable)	Article	Summary
C₁₅H₈O₅	view	view
C₁₄H₂₁N₃O₂S	view	view
C₁₉H₂₁NOS	view	view
C₁₆H₁₅F₂NO₃	view	view
C₁₈H₁₆N₂O₃	view	view
C₁₇H₂₈O₅	view	view
C₉H₁₁FN₂O₃	view	view
C₉H₈N₂O₂	view	view
C₁₂H₁₃NO₅	view	view
C₁₆H₁₈N₄O₇S	view	view
C₁₉H₁₇ClN₄O₂S	view	view

<<- <- -> ->> Prev 1/1 Next



It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

Is CrystalEye a Repository?

Browse

View

1,3-Dihydroxy-9,10-dioxo-9,10-dihydroanthracene-2-carbaldehyd

[OPEN DATA](#)

[<< Table of Contents](#)

Publisher: Acta Crystallographica
Journal: Section E
Year/Issue: 2008/03-00

Article (via DOI): [10.1107/S1600536808004169](https://doi.org/10.1107/S1600536808004169)
Compound Class: organic
Date Recorded: 2008-01-28

Contact Author: Retailleau, Pascal
e-mail: pascal.retailleau@icsn.cnrs-gif.fr

Data collection parameters

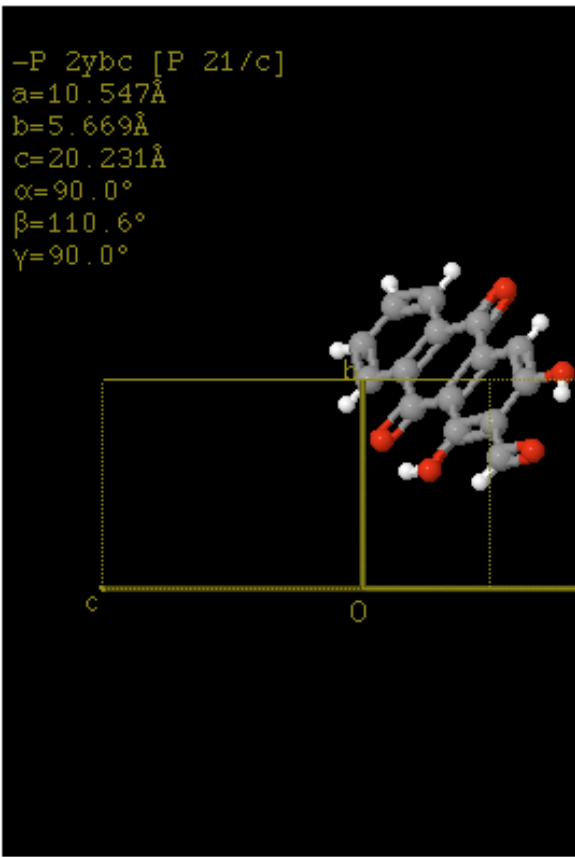
Chemical formula sum	C ₁₅ H ₈ O ₅
Chemical formula moiety	C ₁₅ H ₈ O ₅
Crystal system	monoclinic
Space group H-M	P 21/c
Space group Hall	-P 2ybc
Data collection temperature	293.0

Refinement results

R Factor (Obs)	0.051
R Factor (All)	0.081
Weighted R Factor (Obs)	0.132
Weighted R Factor (All)	0.15

Unit cell parameters:

```
-P 2ybc [P 21/c]
a=10.547Å
b=5.669Å
c=20.231Å
α=90.0°
β=110.6°
γ=90.0°
```



Show no. of unit cells along axis:

a:
b:
c:

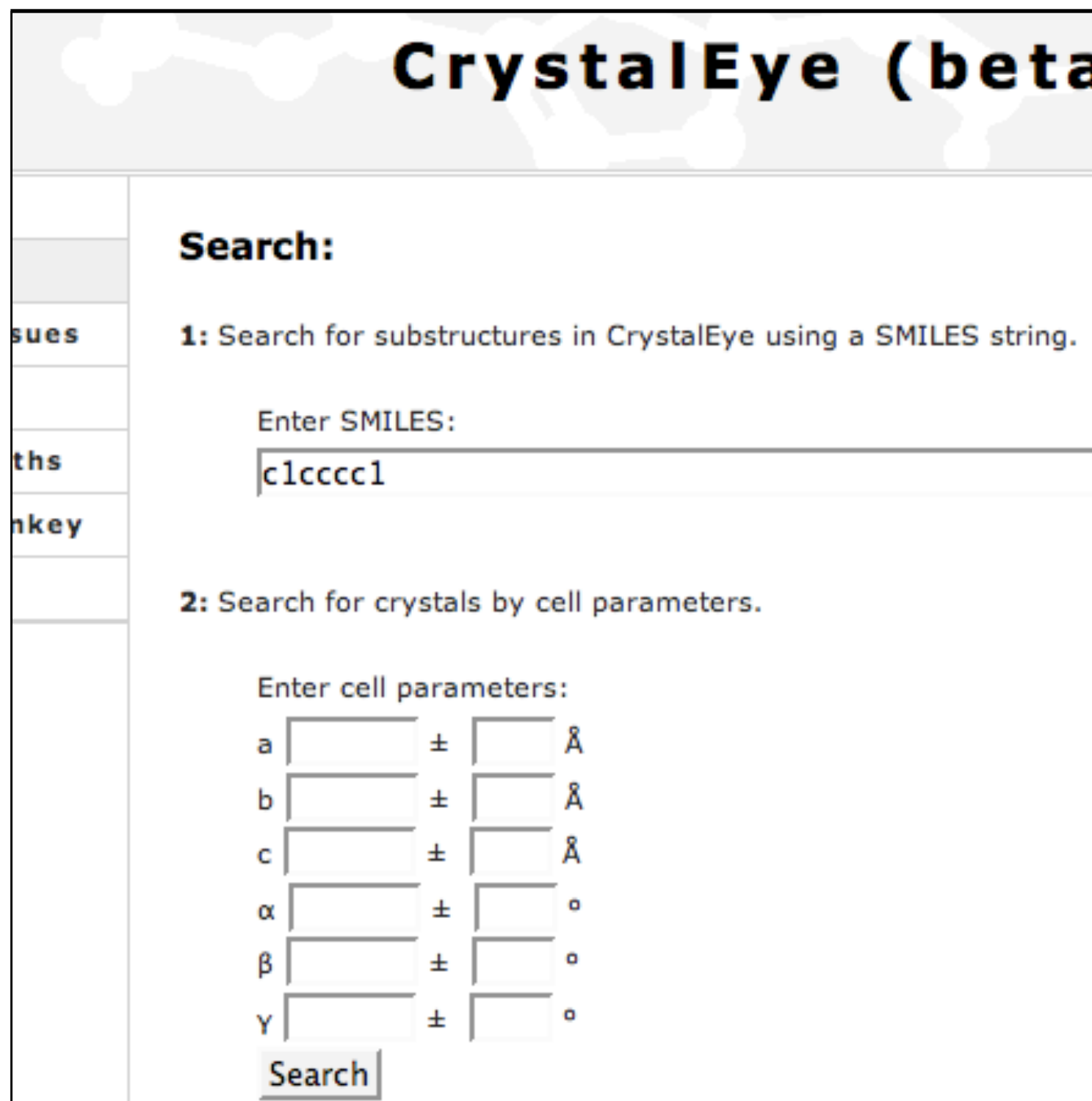
Enter Jmol script:

```
load./bg2162supl_I.complete.cml.x
```

It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

Is CrystalEye a Repository?



CrystalEye (beta)

Search:

1: Search for substructures in CrystalEye using a SMILES string.

Enter SMILES:

2: Search for crystals by cell parameters.

Enter cell parameters:

a	<input type="text"/>	±	<input type="text"/>	Å
b	<input type="text"/>	±	<input type="text"/>	Å
c	<input type="text"/>	±	<input type="text"/>	Å
α	<input type="text"/>	±	<input type="text"/>	°
β	<input type="text"/>	±	<input type="text"/>	°
γ	<input type="text"/>	±	<input type="text"/>	°

Browse

View

Search

It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

Is CrystalEye a Repository?

The screenshot displays the CrystalEye web interface. On the left is a sidebar with navigation links such as 'Items', 'Bradley', 'Discover', and various category links like 'son's weblog (2)', 'entral (1)', 'i - Recent ... (2)', '(1)', '(16)', 'lar (2)', 'All Struc...', 'Space (2)', 'e's Recent ... (1)', 'nouncements (8)', 'cation Today (2)', 'ts.com (5)', 'g - (5)', and 'n (17)'. The main content area is titled 'CrystalEye: All Structures' and includes a search bar, a dropdown menu set to 'All items', and buttons for 'Show: 0 new items - all items', 'Mark all as read', and 'Refresh'. Below this is a list of items, with the first item selected and its details shown. The details include a star icon, a title 'Summary page for crystal structure from DataBlock b719624esup1 from article b719624e in issue 2008 Society of Chemistry, Organic and Biomolecular C', a chemical structure diagram, and a link to the 'Original enclosure'. The chemical structure is a complex polycyclic molecule with a benzene ring fused to a five-membered ring, which is further fused to a six-membered ring containing a nitrogen atom (HN) and an oxygen atom (O). There are also two phenyl rings attached to the structure. The interface also features a sidebar on the right with links like 'Add star', 'Share', 'Email', 'Mark as read', and 'Edit tags: d'.

Browse

View

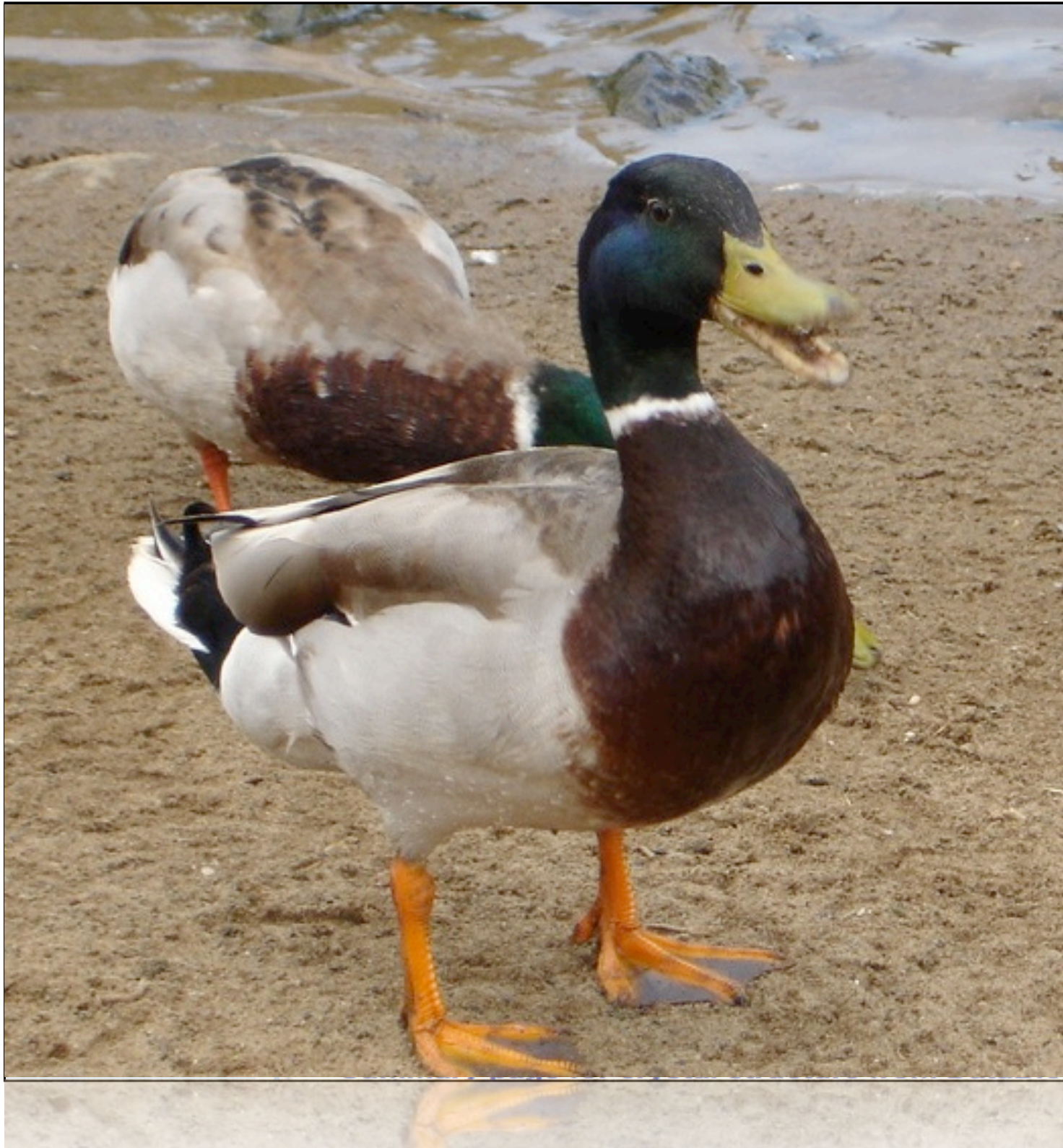
Search

Harvest

It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

Is CrystalEye a Repository?



Browse

View

Search

Harvest

It walks like one, quacks like one, swims like one and goes well with orange sauce like one...

CrystalEye is a “just” a web site, but has many of the features you’d expect from a repository.

CrystalEye and the Subject Repository

e-Crystals Federation

Simon Coles has presented on e-Crystals already this morning. Over the course of e-Crystals we'll be working out how CrystalEye can work as part of a federated pan-institutional subject repository.

CrystalEye and the Institutional Repository

How could / should the University of Cambridge institutional repository interact with CrystalEye?

CrystalEye is a data repository, or at the very least an amply-featured data collection.

The way it's been most usually done in the past has been to move or copy the data over to a centralized repository. We don't believe this would work for most long-tail science applications; there's too much domain expertise involved in curating the data, and it would make managing the data more difficult.

Clifford Lynch's original definition of the Institutional Repository:

IR Services: Some Ideas

Technological

- Additional backup
- Preservation
- URL redirection and management
- Promote access through portals
- Equipment

Organisational

- Create data management roles in departments
- Data management training for academics
- “Transferrable skills” training for PhD students?

These are just ideas from our perspective.

I'd like to emphasize how useful data training could be – Nick lost plenty of time due to data management mistakes, particular around identifiers and discarding intermediate data.



Conclusions

26

Short version:

- * Keep it simple
- * Share your data to let the community drive functionality
- * Focus on the web

Long version:

So why did CrystalEye teach me to love the web?

Because we never had to stop and think “now we have to stop and get this data into a

Copyright and Notices

Flickr images:

PMR: borazivkovic (BY)

Crystal: wabberjocky (BY-NC-SA)

Spider: dincordera (BY-NC-SA)

Car: snellgrove (BY-NC-SA)

Ascent of man : Kaptain Kobold
(BY-NC-SA)

Open/Close: mag3737 (BY-NC-SA)

Spices: estherase (BY-NC-SA)

Combine harvester: tricky TM (BY-NC-SA)

Sunset: nexus6 (BY-NC-SA)

Horse race: Chris Breeze (BY)

Duck: Frogzone 1 (BY-NC)

Web/Pylon: Orcoo (BY-NC-SA)

Others:

Atom icon: Mozilla Foundation
(MPL)

Credits and Links

- Thanks to Nick Day and Peter Murray-Rust for all their work on CrystalEye
- Thanks to Andrew Walkingshaw for his work on the CrystalEye RDF and to Talis for the use of their Platform Store
- CrystalEye: <http://wwmm.ch.cam.ac.uk/crystaleye/>
- Coverage on some of the features mentioned here on my blog: <http://wwmm.ch.cam.ac.uk/blogs/downing/>
- Coverage on CrystalEye and much, much more on PM-R's blog: <http://wwmm.ch.cam.ac.uk/blogs/murrayrust/>