# Semantic Chemical Publishing

**Nick Day\*, Peter Corbett, Peter Murray-Rust**
Unilever Centre for Molecular Informatics, University of Cambridge, UK.
March 27th, 2007

- All software Open Source
- \* ned24@cam.ac.uk

# Overview

- What is 'semantic chemistry' and markup?

- **OSCAR3** – robotic analysis of chemistry in free text
    - recognition of chemical names
    - name-2-structure
    - chemical verbs, adjectives and reaction names
    - terminologies (e.g. techniques)
    - RSC **Project Prospect** …

- **CrystalEye** – creating semantic chemistry from crystallography:
    - High-throughput robotic harvesting
    - Re-use using CIF2CML
    - Dissemination through CMLRSS

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# The Semantic Web

"People keep asking what Web 3.0 is. I think maybe when you've got an overlay of scalable vector graphics […] on Web 2.0 and access to a semantic Web integrated across a huge space of data, you'll have access to an unbelievable data resource."

- Tim Berners-Lee, A 'more revolutionary' Web (2006)

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# …Let's change the vision to chemistry…

# The Chemical Semantic Web

"…when you've got an overlay of

…..

- on Web 2.0 and access to a semantic Web integrated across a huge space of data, you'll have access to an unbelievable data resource."

*[our  adaptations]*

*… what are chemical semantics and CML?…*

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# Implicit and explicit semantics

- *Implicit semantics*

  "**Compound 2a melted** at **119ºC**"

  *humans are good at interpreting this; machines see just a string.*

- *Explicit semantics*

  **CML Schema**

  ```
  <cml:molecule ref="2a">
    <cml:property>
      <cml:scalar dictRef="prop:mpt"
          units="units:celsius"
          dataType="xsd:float"
      >119</cml:scalar>
    </cml:property>
  </cml:molecule>
  ```

  **Molecules in CML/InChI**

  **propertyDictionary**

  **unitsDictionary**

  **W3CSchema**

  *4 namespaces, 3 dictionaries*

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# UCC's approach to creating Semantic Chemistry

- Authoring tools for theses and collaboration with publishers
- XML-ization (through FoX) of Comp. Chem. codes (MOPAC, CASTEP, SIESTA, GULP, ABINIT, DL_POLY GAMESS…)
- Capturing/conversion of CML data at source (SPECTRa)
- Rich clients (Bioclipse)
- Legacy Conversion (OpenBabel, CDK, JUMBO…)
- Intelligent Ontologies (Golem)

*(today we will cover the following…)*
- **Chemical Linguistics and text-mining (OSCAR3)**
- **Legacy Conversion – CIF**

*Many of these semantic chemical components are now deployed or prototyped…*

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# Chemical semantic framework at UCC

# "OSCAR"

2-Chloro-5,11-dihydro-11-ethyl-8-chloromethyl-6H-dipyrido[3,2-b:2',3'-e][1,4]diazepin-6-one (**15**)

A suspension of **14** (177 mg, 0.58 mmol) in $CH_2Cl_2$ (100 mL) was treated with thionyl chloride (0.3 mL) followed by triethylamine (1 mL). The reaction mixture was stirred at room temperature for 1 ... tion was obtained. Then, saturated aqueous $NaHCO_3$ was added and the mixture ... $CH_2Cl_2$. The organic layer was washed with water, then dried ($Na_2SO_4$) and ... educed pressure. The residue was purified by silica gel column chromatography (20%EtOAc/hexane) to give **15** (163 mg, 87%) as a pale yellow solid, m.p. 226-227 °C; FTIR (KBr)

- Recognition of chemical entities.
- Name2structure, chemical diagrams, canonical identifiers
- Chemical heuristics to parse article full-text
- Links to ontologies and molecular databases.
- Open source
- High-throughput – 500, 000 PubMed abstracts parsed
- Substructure and similarity search on corpora

OSCAR1 + CheckCML (2003, 2004, 2005, 2006) Student projects supported by RSC
SciBorg (2005-2009)  EPSRC project (Computer Lab, Chemistry, Cambridge)
**OSCAR3 (Peter Corbett**)

**UNIVERSITY OF CAMBRIDGE**

**Unilever Cambridge**
Centre for Molecular Informatics

# OSCAR3 Concepts - Example



All markup is **automatic**

**…how can this be used for publishing?...**

**UCC and RSC have been collaborating on transferring this technology to journal articles…**
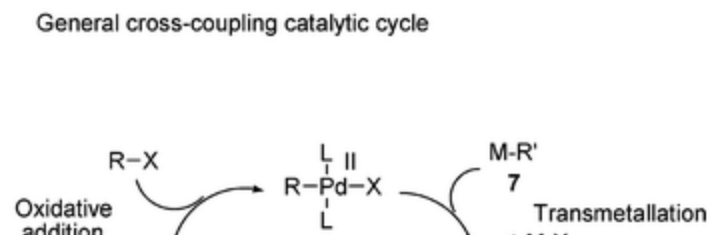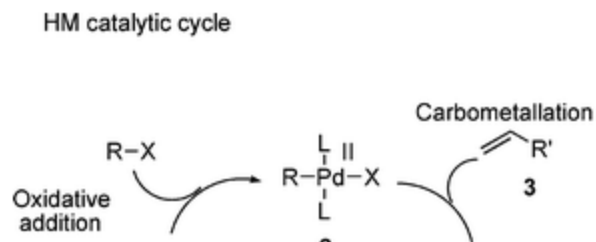
Project Prospect (2007) adds semantics…

# Project Prospect (RSC)

Typical HTML paper                                    Semantics confined to hyperlinks

There has been a broadly accepted understanding of the mechanism operating in the HM reaction for many years; [http://www.rsc.org/delivery/_ArticleLinking/ArticleLinking.asp?JournalCode=OB&Year=2007&ManuscriptID=b611547k&Is...] (Show/Hide) Toolbox initial oxidative addition of the halide to a palladium(0) catalyst. Despite various claims for a possible palladium(II/I)...

the evidence for this is poor, since it has been shown that in the majority of cases, the palladacycles involved act as reservoirs of palladium some of which is reduced to palladium(0). Further evidence against this mechanism comes from gas phase computational studies which ind... the rate determining step in a palladium(II/IV) cycle involving iodobenzene would be the oxidative addition of iodobenzene to palladium.[45]

the actual rate determining step in the HM reaction of aryl iodides is not oxidative addition [46] (*vide infra*) this indicates that a palladium(II... is not in operation. Hence, the mechanism of the HM process can be represented by Scheme 1, involving a palladium(0) species 1 underg... oxidative addition to generate a palladium(II) species 2, which reacts with the olefin component 3, possibly following initial $\eta^2$-coordinatio... palladium atom. This results in a carbometallation reaction to generate palladium(II) alkyl complex 4. Elimination of palladium hydride from 4 furnishes the product 5 and base assisted elimination of HX from palladium(II) complex 6 regenerates the active palladium(0) catalyst 1.

HM catalytic cycle

General cross-coupling catalytic cycle

Carbometallation

R–X    L II    R'
     R–Pd–X       3
Oxidative    L
addition

R–X    L II    M–R'
     R–Pd–X       7
Oxidative    L    Transmetallation
addition

**Prospect adds more semantics…**

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

## … **Prospect markup includes:**

- CML (Chemical Markup Language)

- InChI

- IUPAC Gold Book

- Gene Ontology

- …and more coming …

…the marked-up semantic paper…

# Project Prospect RSC



… *but not all chemistry is in free text …*

# Chemistry is also "Data"

### 10.6 Strukturdaten auf B3LYP/6-31G(d) Niveau

2-Pyridon (1a)

...\C5H5N1O1\HARRY\12-Jul-2000\0\\#BECKE3

```
loop_
 _atom_site_label
 _atom_site_type_
 _atom_site_fract
 _atom_site_fract
 _atom_site_fract
 _atom_site_U_iso
 _atom_site_adp_t
 _atom_site_occup
 _atom_site_symme
 _atom_site_calc_
 _atom_site_refir
 _atom_site_disor
 _atom_site_disor
C1 C 0.2103(4) 0.
H1 H 0.2003 0.914
C2 C 0.0615(5) 0.
H2 H -0.0462 0.84
```
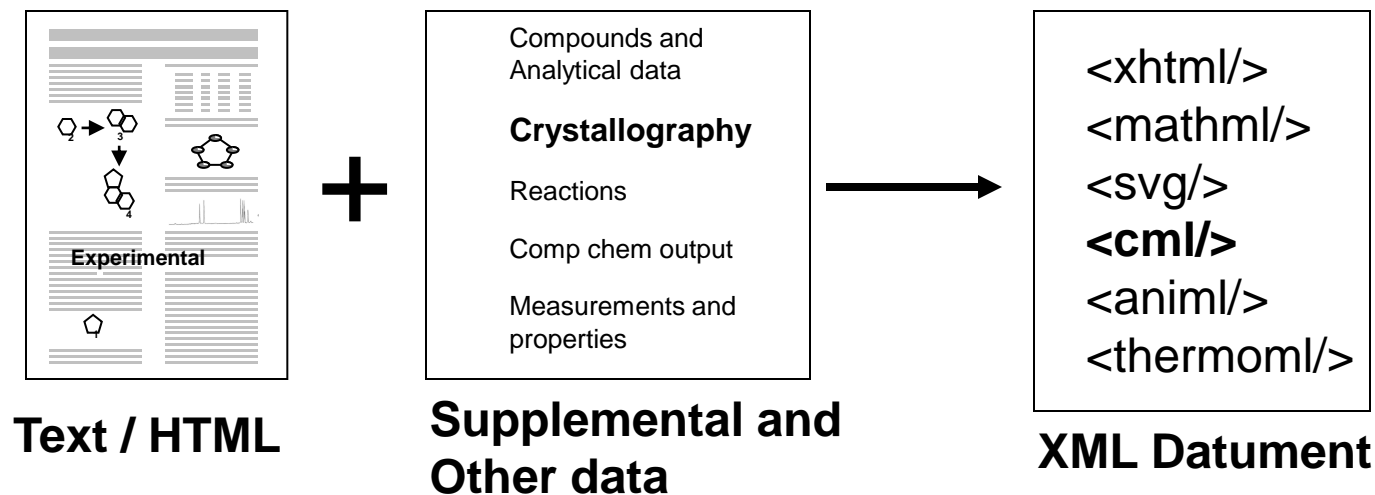
Tabelle 10.1: Beobachtete Geschwindigkeitskonstente für die Reaktion von Butylamin (3) (250 mmol/l) mit p-Nitrophenylacetat (2) (50 mmol/l) in Deuterochloroform bei 23 °C, katalysiert mit den 2-Pyridonen 1, 7, 9, und 10.

**… some examples of data taken from theses.**

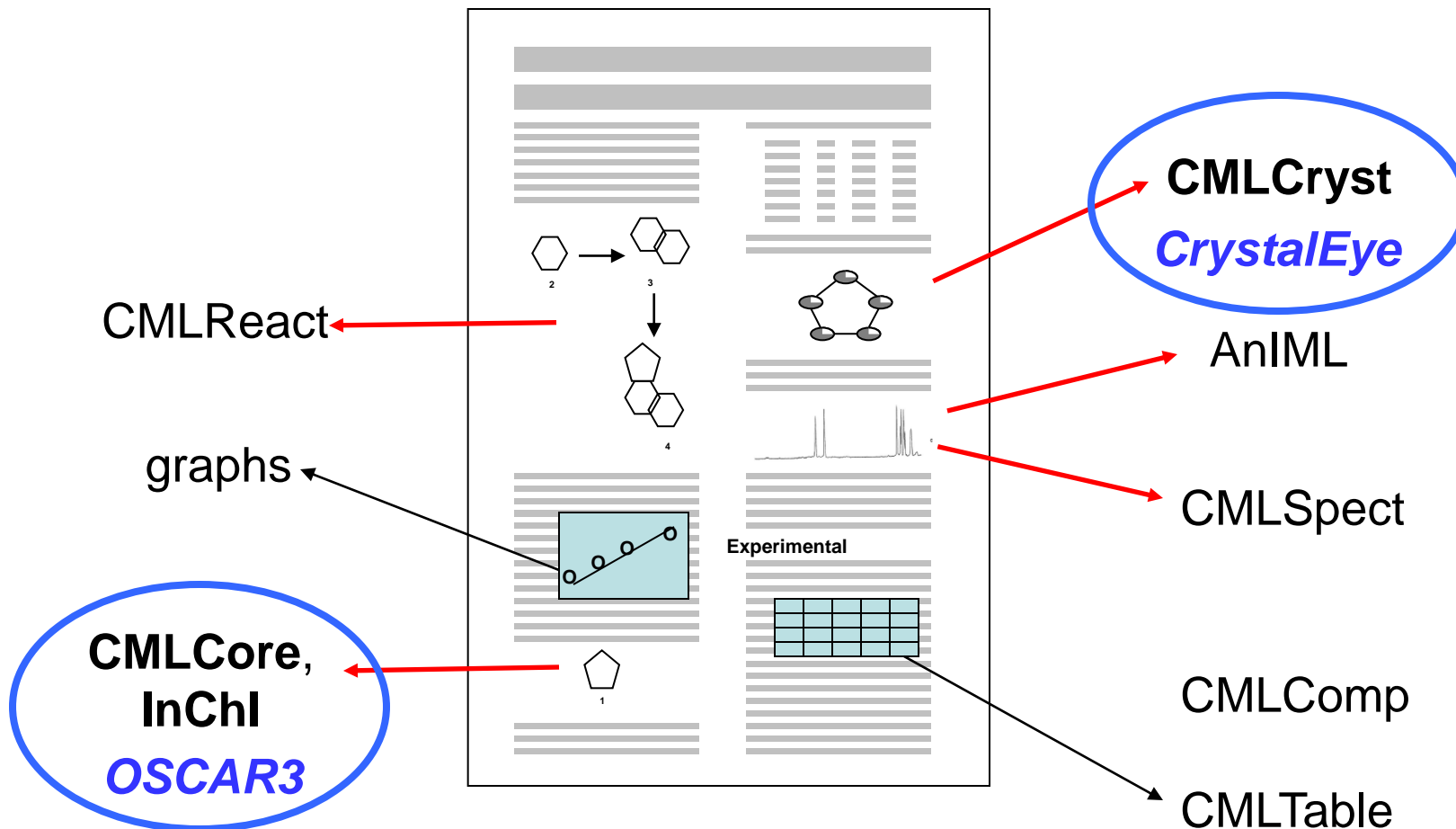| [1] | | | | | | [10] | $k_{obs}$ |
|---|---|---|---|---|---|---|---|
| [mmol/l] | [$10^{-4}$ s$^{-1}$] | [mmol/l] | [$10^{-4}$ s$^{-1}$] | [mmol/l] | [$10^{-4}$ s$^{-1}$] | [mmol/l] | [$10^{-4}$ s$^{-1}$] |
| 0.34 | 7.30 | 0.200 | 6.97 | 0.500 | 6.90 | 0.063 | 6.34 |
| 0.67 | 7.31 | 0.400 | 6.99 | 1.25 | 6.94 | 0.10 | 6.54 |
| 1.66 | 8.42 | 0.500 | 7.14 | 2.50 | 7.46 | 0.25 | 6.75 |
| 3.00 | 10.42 | 1.00 | 7.13 | 3.75 | 7.70 | 0.38 | 8.12 |
| 3.30 | 10.21 | 2.50 | 7.63 | 10.0 | 8.50 | 0.50 | 8.45 |
| 5.00 | 11.41 | 4.00 | 8.40 | 15.0 | 8.76 | 1.00 | 7.71 |
| 16.0 | 16.08 | 5.56 | 8.56 | 33.33 | 9.21 | 1.25 | 8.94 |

# Semantic Chemistry and the Datument

The **document** is only part of the scientific record
We can transform the experimental **data** to CML.
The integrated result is a **datument**…



**Text / HTML**                 **Supplemental and Other data**                 **XML Datument**

- P. Murray-Rust, *The complete chemical E-publication*, 216th ACS National Meeting, Boston, August 23-27 (**1998**), CINF-033.
- Peter Murray-Rust, Henry S. Rzepa and Michael Wright, Development of Chemical Markup Language (CML) as a System for Handling Complex Chemical Content, *New J. Chem.*, **2001**, 618-634.
- P. Murray-Rust and H. S. Rzepa, "The Next Big Thing: From Hypermedia to Datuments", *J. Digital Inf.*, **2004**, **5**, article 248, 2004-03-18.

# The Datument



CMLCryst
*CrystalEye*

CMLReact

AnIML

graphs

CMLSpect

CMLCore, InChI
*OSCAR3*

Experimental

CMLComp

CMLTable

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# Chemical Crystallography

Universally published as **CIF**.

- Complete output of structure experiment,
  ~~~ary~~~ ary data for article full-text.
  ~~~hed~~~ hed online per year
  ~~~d~~~ d per year (e.g. theses)

```
data_abU8

_chemical_formula_sum            'C13 H17 N O4 S'
_chemical_formula_weight         283.34

_symmetry_space_group_name_H-M   P2(1)/a

loop_
  _symmetry_equiv_pos_as_xyz
  'x, y, z'
  '-x+1/2, y+1/2, -z'
  '-x, -y, -z'
  'x-1/2, -y-1/2, z'

_cell_length_a                   8.367(4)
_cell_length_b                   19.764(8)
_cell_length_c                   8.672(4)
_cell_angle_alpha                90.00
_cell_angle_beta                 95.16(3)
_cell_angle_gamma                90.00
_cell_volume                     1428.2(11)
loop_
  _atom_site_label
  _atom_site_type_symbol
  _atom_site_fract_x
  _atom_site_fract_y
  _atom_site_fract_z
  _atom_site_U_iso_or_equiv
  _atom_site_thermal_displace_type
  _atom_site_occupancy
  _atom_site_calc_flag
  _atom_site_refinement_flags
  _atom_site_disorder_group
N1 N 0.0887(2) 0.55103(9) 0.6991(2) 0.0399(5) Uani 1 d . .
H1 H 0.1327(2) 0.59208(9) 0.6504(2) 0.067(8) Uiso 1 calc R .
```

```xml
<entry dataType="xsd:double" minInclusive="0.0"
  maxInclusive="180.0" units="units:deg
  unitType="unitType:angle" id="_cell_angle_alpha">
  <definition>Unit-cell angles of the reported
    structure in degrees. The values of
    _refln_index_h, *_k, *_l must correspond to the
    cell defined by these values and _cell_length_a,
    *_b and *_c. The values of
    _diffrn_refln_index_h, *_k, *_l may not
    correspond to these values if a cell
    transformation took place following the
    measurement of the diffraction intensities. See
    also _diffrn_reflns_transf_matrix_.</definition>
  <scalar dictRef="iucr:category">cell</scalar>
</entry>
```
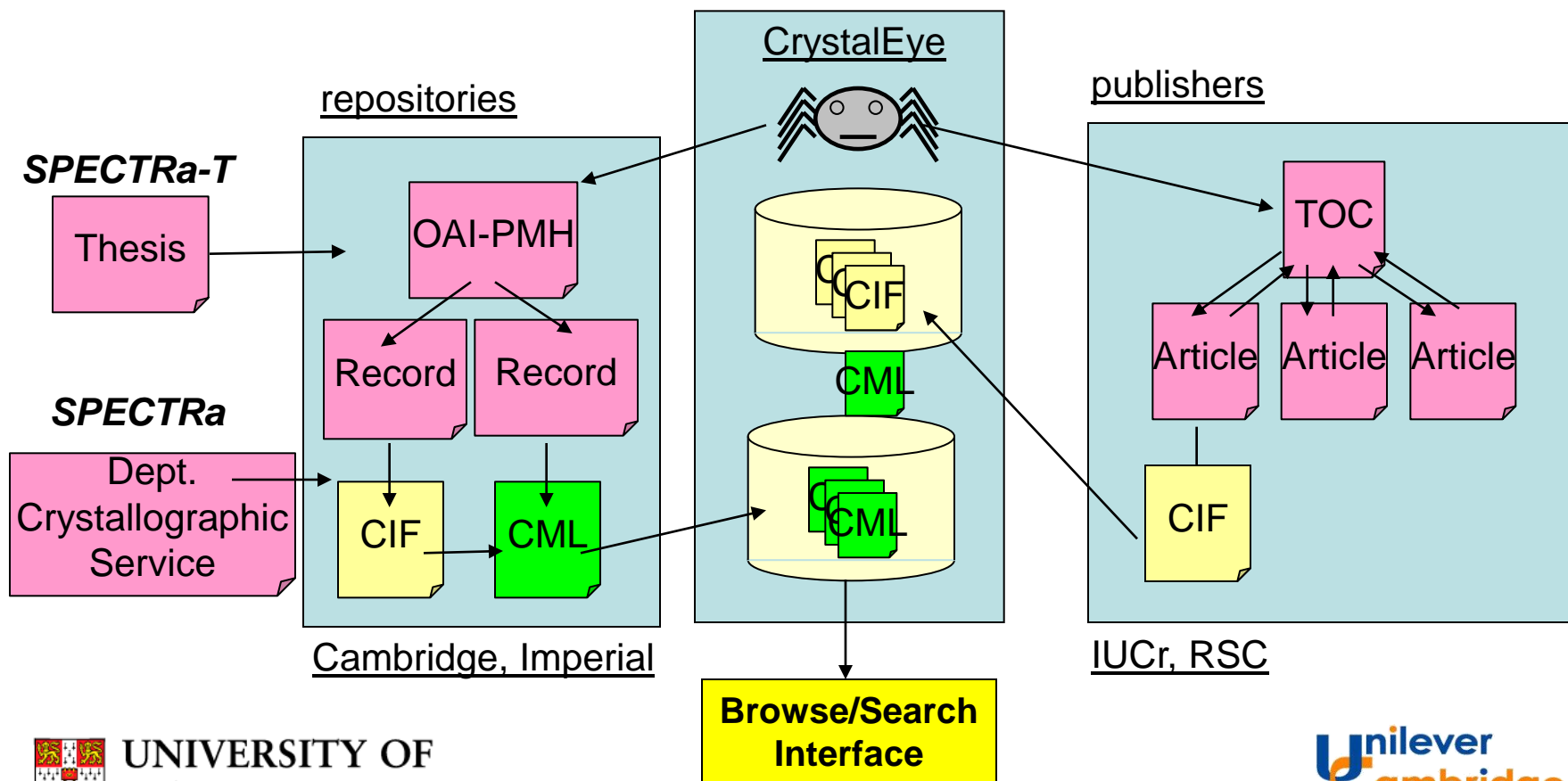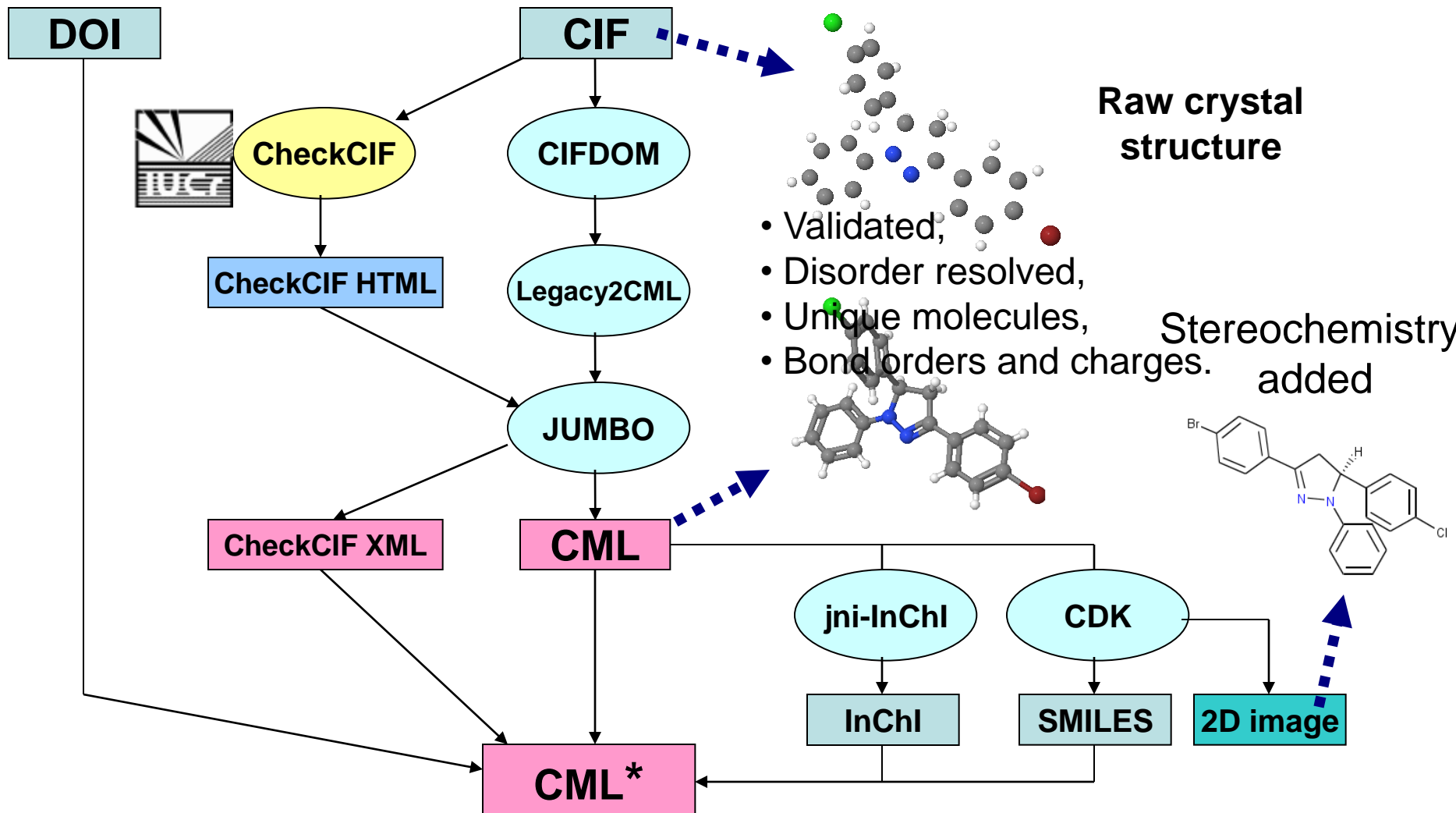
# CrystalEye

**The aim:**

To automatically create semantic chemistry from crystallography (CIFs) published on the Web.

# Aggregation

- Web spider checks publishers and repositories every day.
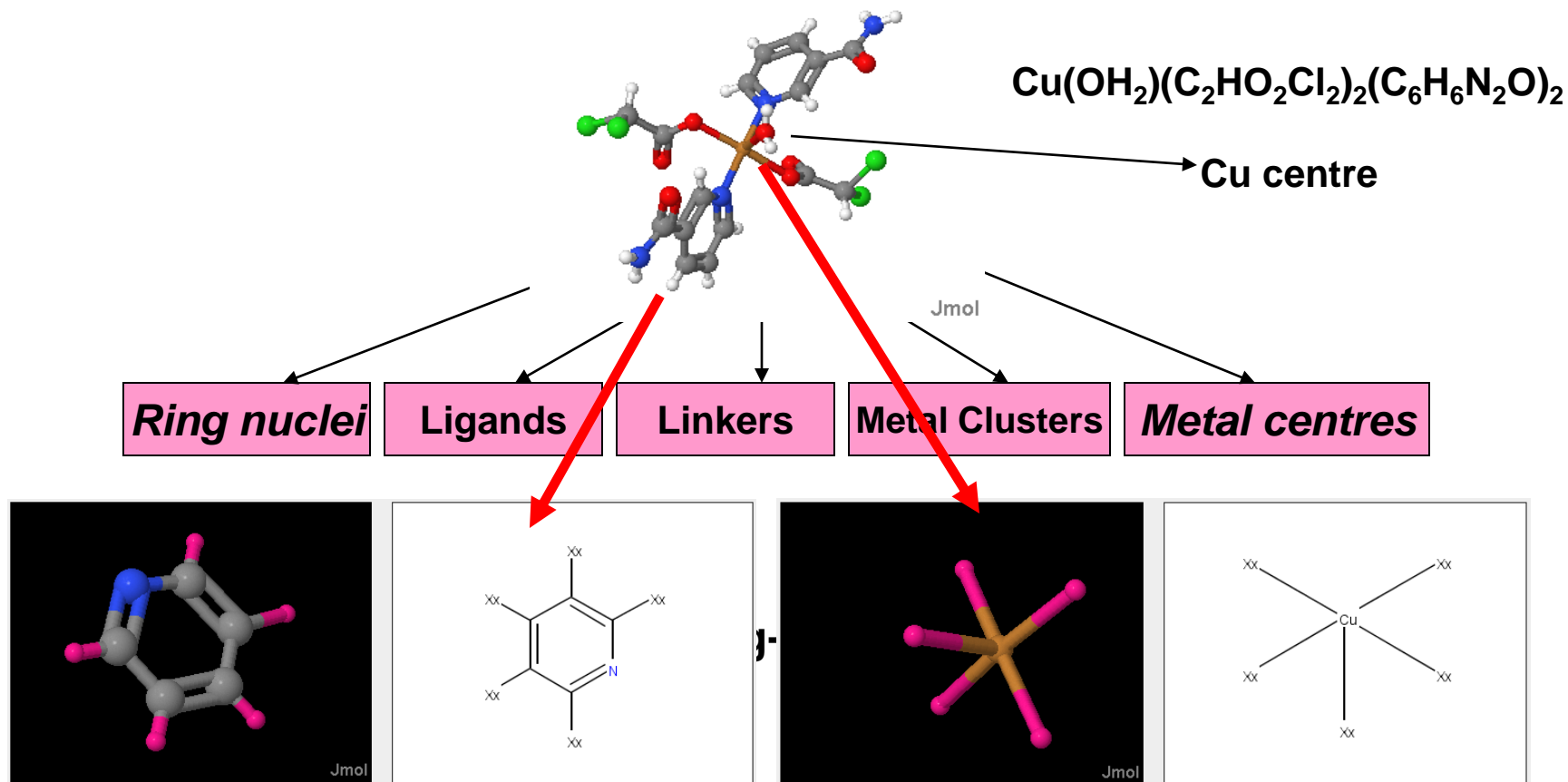- Currently over **60,000 validated** CIF files.
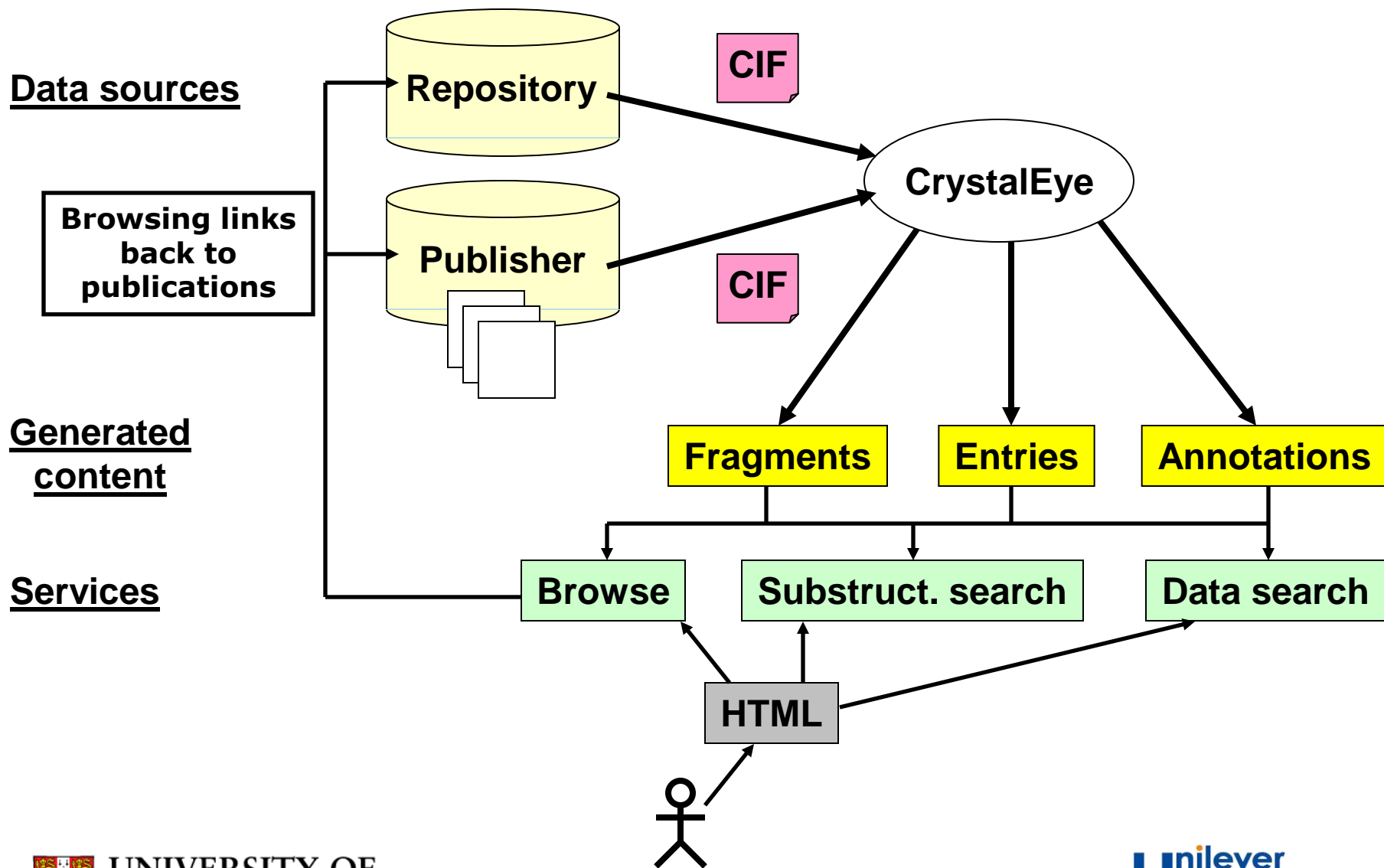
# Marking up and Validation



**DOI**

**CIF**

**CheckCIF**

**CheckCIF HTML**

**CIFDOM**

**Legacy2CML**

**JUMBO**

**CheckCIF XML**

**CML**

Raw crystal structure

- Validated,
- Disorder resolved,
- Unique molecules,
- Bond orders and charges.

Stereochemistry added

**jni-InChI**

**CDK**

**InChI**

**SMILES**

**2D image**

**CML***

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# CrystalEye: Re-Use through XML/CML

- Automatic generation of fragments

$Cu(OH_2)(C_2HO_2Cl_2)_2(C_6H_6N_2O)_2$

**Cu centre**

Jmol

| *Ring nuclei* | Ligands | Linkers | Metal Clusters | *Metal centres* |
|---|---|---|---|---|



- ca. 1 million fragments with 50,000 different chemical types
- Open Access via automatically generated HTML

# CrystalEye for humans

# CrystalEye Webpage Demo

Let's assume we're interested in Cu-N bonds:

- Browse to title page

- View structure data

- Explore fragments

- Inspect bond lengths



UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# CMLRSS newsfeeds

- How can the chemist find every Cu-N bond immediately?

Structure:

- = contains Cu-N
- = no Cu-N

*read*

**CrystalEye**

**TOC**   **Cu-N feed**

...

**Web browsing**

**What we've just been doing**

**Using RSS**

RSS Reader

- CrystalEye uses RSS 1, RSS 2 and Atom 1.0 to create both RSS and CMLRSS feeds.

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# CrystalEye: KnowledgeBase not DataBase

- Aggregation by robots, not humans
- All types of chemistry (organic, inorganic, etc…)
- Social computing - aggregates **COD\*,** theses
- Software validation by robots, not humans,
- Open and free;

- Goes live in April…

**\* Crystallographic Open Database**

**UNIVERSITY OF CAMBRIDGE**

**Unilever Cambridge**
Centre for Molecular Informatics

# Chemical semantic framework at UCC

# Cambridge Semantic Chemistry

- Released:
    - CML – XML for chemistry
    - JUMBO – library for CML
    - OSCAR1/CheckCML – data validation
    - OSCAR3 – text mining
    - FoX – Fortran XML

- Soon:
    - CrystalEye – crystallographic knowledgebase
    - SPECTRa – chemical repositories

- Later…
    - Golem - ontologies
    - CMLUnits

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics

# Acknowledgements

- CML - Henry Rzepa
- OSCAR1 – Sam Adams, Joe Townsend, Fraser Norton, Justin Davies, Richard Marsh, Jonathan Goodman
- OSCAR 3 – Ann Copestake, SimoneTeufel (Computer Lab)
- SPECTRa – Jim Downing, Alan Tonge, Peter Morgan
- Software - CDK, Jmol, jni-InChI and many Blue Obelisk contributions
- Timo Hannay (NPG), Richard Kidd (RSC), Colin Batchelor (RSC), Brian McMahon (IUCr)

# Thankyou.

UNIVERSITY OF CAMBRIDGE

Unilever Cambridge
Centre for Molecular Informatics