



# Data quality and the value of structural databases

Colin Groom

Crystallographic Information and Data Management  
A Satellite Symposium to the 28th European Crystallographic Meeting IV.  
Towards ever better science

Sunday, August 25, 2013



# Contents

- The CCDC
- The CSD
- Creating the CSD
- Curating the CSD
- Current challenges
  - Beyond structural chemists
  - Structure publishing
  - Non-semantic data
- Unanswered questions



# The Cambridge Crystallographic Data Centre

- A not-for-profit, charitable institution, established 1965
- Self-financing and self-administering
  - Funded by contributions from user community
- A University of Cambridge Partner Institute
- A UK Government Independent Research Organisation
- Around 50 members of staff
  - Scientists, software developers, IT, finance and user services

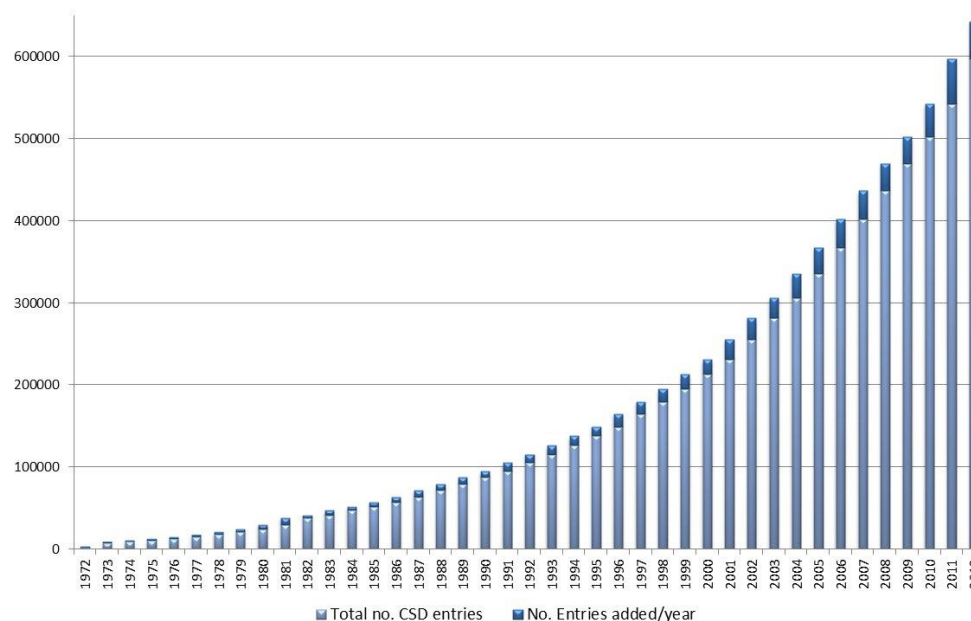
“advancement and promotion of the  
science of chemistry and crystallography  
for the public benefit”





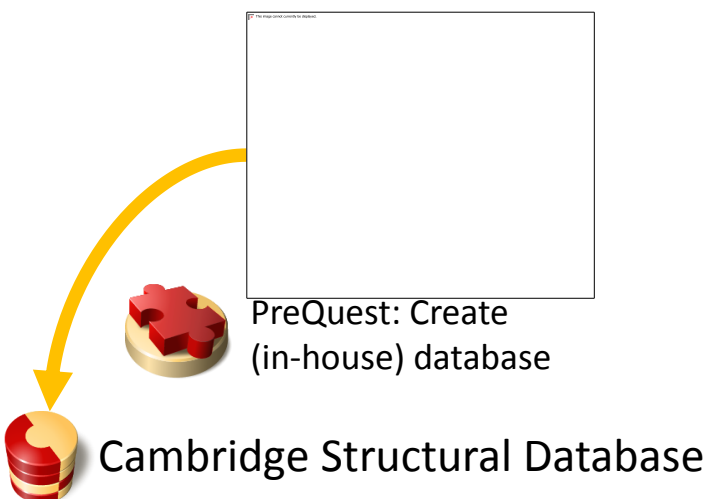
# The Cambridge Structural Database


- All published small molecule crystal structures
  - Currently 658,047
  - Comprehensive
- All curated
  - Corrected
  - Standardised
  - Enriched
- All searchable
  - ‘Chemistry’ added
  - Sophisticated tools

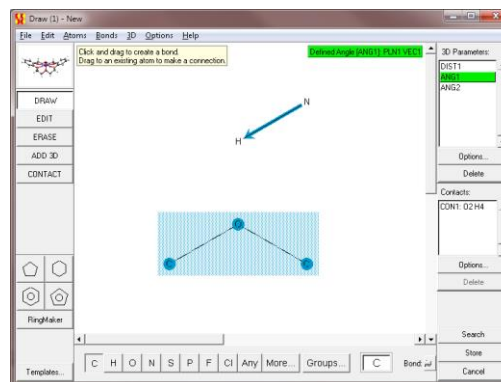





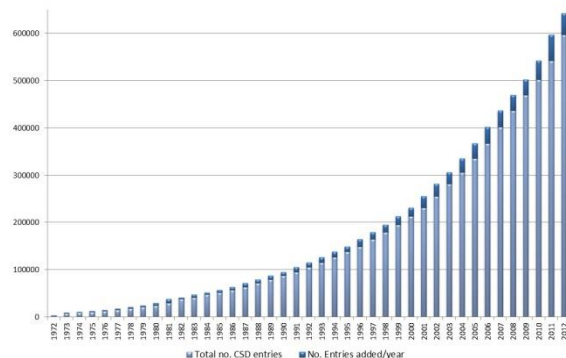
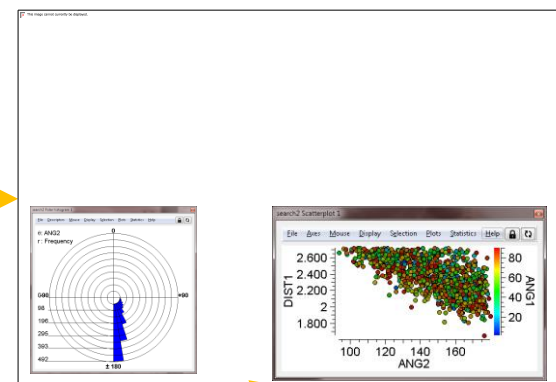
# The Cambridge Structural Database System




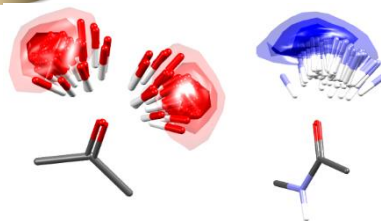
 ConQuest: Advanced 3D searching




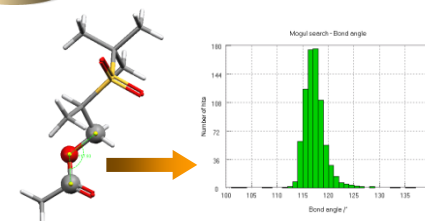
 Mercury: Visualisation & data analysis



 IsoStar: Molecular interaction analysis



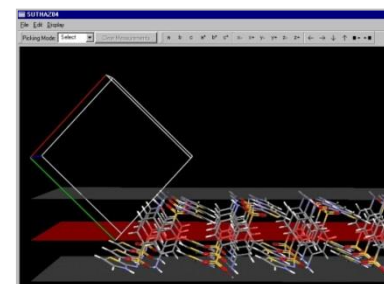
 Mogul: Molecular geometry analysis





# Services provided

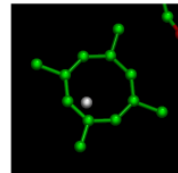
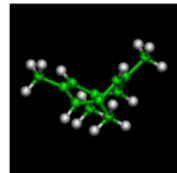
- Access to original deposited data
  - Individual data sets available to anyone
  - No cost archiving and dissemination of crystallographers output
    - c.f. Publication APCs
- Software tools
  - enCIFer (validation of CIFs)
  - Mercury (crystal structure visualisation)
  - Ligand dictionaries
- Targeted subsets of curated data
  - Teaching subsets
  - PDB subset
- CSD System



AROMATICITY > STEPS REQUIRED > Consider what happens when we treat cyclooctatetraene with a powerful reducing agent.

**Consider what happens when we treat cyclooctatetraene with a powerful reducing agent**

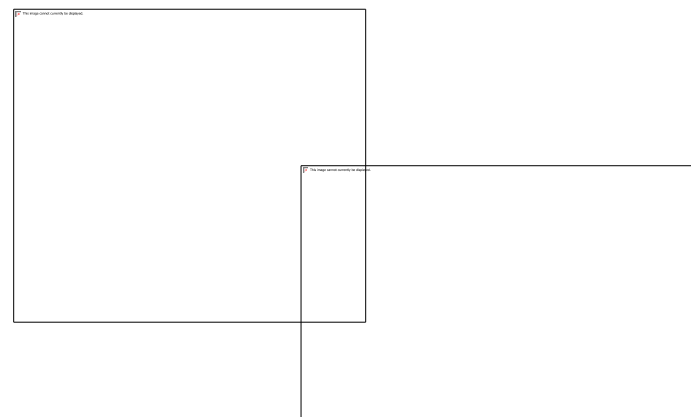
- If 1,3,5,7-tetramethylcyclooctatetraene (refcode TMCOTT) is treated with alkali metals a dianion is formed (refcode TMCCKE).
- Look closely at the structures of 1,3,5,7-tetramethylcyclooctatetraene (refcode TMCOTT) and the resultant dianion (refcode TMCCKE). How do these two compounds differ structurally?
- You should find that the dianion is planar and all bonds lengths are equivalent (within experimental error). Whereas the neutral compound is non-planar ("tub" shaped) with alternate double and single bonds lengths of 1.48Å and 1.33Å.

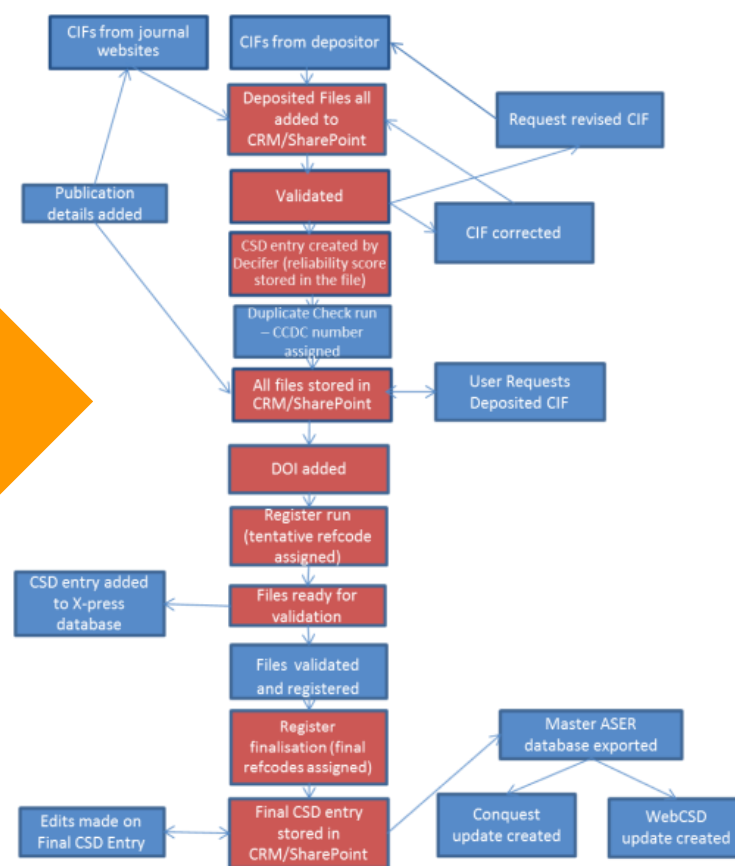


left: "tub" shape of 1,3,5,7-tetramethylcyclooctatetraene (refcode TMCOTT), right: the resulting planar dianion (refcode TMCCKE)



- Structures from depositors and journals
  - Referees, embargoes, revision
- Curation
  - Coordinates to ‘chemistry’
  - Duplication
  - Enrichment
- Distribution and access
  - Searching, analysis, application









# CSD-Xpedit

Microsoft Dynamics CRM

File Edit View Favorites Tools Help

Convert Select

Records Collaborate Process Data ReMAP

Workplace

My Work

- Dashboards
- Activities
- Calendar
- Duplicate Detection
- Queues

Workplace

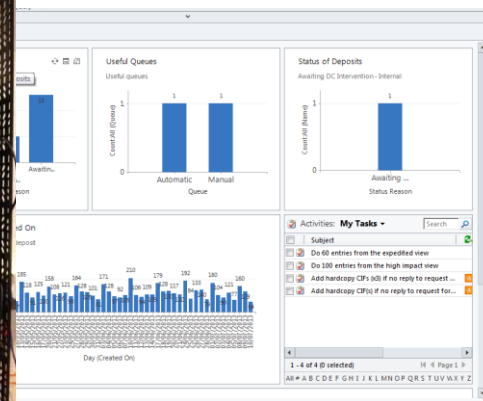
- Sales
- Marketing
- Service
- Settings
- Resource Center

Deposits: Awaiting DC Intervention - External

Name	Status Reason	Depositor	Deposition Co...	Subject (Email)
1283576-DEP	Awaiting CCDC Numbers	Shun-Yi Wang	CCDC Deposit	New, valid web deposit
1283574-DEP	Awaiting Correction	Konstantin Doma	CCDC Deposit	Re: CCDC Depository Req
1283565-DEP	Awaiting CCDC Numbers	Michael Jennings	CCDC Deposit	New, valid web deposit
1283562-DEP	Awaiting CCDC Numbers	Laura Cañadillas-I	CCDC Deposit	pre-publication cif file
1283561-DEP	Awaiting Correction	Dra. Victoria Cast	CCDC Deposit	Deposited Data - CCDC 9
1283560-DEP	In progress of assigning ...	Mareike Jahnke	CCDC Deposit	Deposit of 7 X-ray structu
1283556-DEP	Awaiting Correction	Chang Hong	CCDC Deposit	Re:CCDC Depository Req
1283526-DEP				
1283440-DEP				

1 - 9 of 9 (0 selected)

All # A B C



Publication

**Eur.J.Inorg.Chem. (2013) ,3316**

Name Eur.J.Inorg.Chem. (2013) ,3316



Curated Data: Curated Data Associated View

Name	Deposit	CCDC Number	Source File (External Format)
1643223-C...	1000218-DEP	929,193	structure_1999264251.txt
1762050-C...	1274796-DEP	929,192	ccdc.cif
1762051-C...	1274796-DEP	929,191	dvshsun11vls_mach.cif
1643224-C...	1000218-DEP	929,192	structure_1999264250.txt
1643225-C...	1000218-DEP	929,191	structure_1999264248.txt

Overview Register All Text Diagram Visualiser All Hits

Chemical structure diagram showing a complex molecule with a central core and various substituents.

CCDC Number	CCDC 910817	Identifier/CCDC Refcode	1660396/REVJAG	Previous Refcode
Space Group	P -1	Z, Z'	Z: 1 Z': 0.5	R-Factor (%)
Temperature (K)	393	Temperature (text)	393 K	Pressure
Formula weight (CCDC)	632.745	Density (CCDC)	1.1895	Cell Volume
Published Formula Weight	632.73	Density (author)	1.189	Published Cell Volume
Powder Study	unknown	Radiation Probe	unknown	Radiation Source
Color	colorless	Habit	Blocks	Melting Point

Literature Reference

T.V.Sreevidya, Deng-Ke Gao, T.Lavay, M.Botoshansky, M.Kafant, *Cryst.Growth Des.* (2013), **13**, 936, doi:10.1021/cg3016707

Compound Name

1,1,6,6-Tetraphenyl-2,4-hexadiene-1,6-diol bis(6-methyl-2(1H)-pyridinone)

Synonym

Formula

C<sub>30</sub>H<sub>22</sub>O<sub>2</sub>Zn(C<sub>6</sub>H<sub>5</sub>N O)

Cell Lengths

a 8.4750(5) b 10.8530(6) c 11.3600(10)

Cell Angles

α 64.798(3) β 81.182(2) γ 69.185(7)

Disorder

C16A,C18A,C20A and C16B,C18B,C20B disordered over two sites with occupancies 0.627/0.373.

Errors

Total of 1 CSD Editor

Warnings

Line 95 CifDatabase

Data value should be a number; diffm\_standards\_number

Cell contents: C30 H22 O2,4(C6 H7 N O1); Formula unit: C30 H22 O2,4(C6 H7 N O1); Z=1, Z'=0.50; perfect match

Rejected Hits

Accepted Hits

Comparison Entries

Incoming Entries



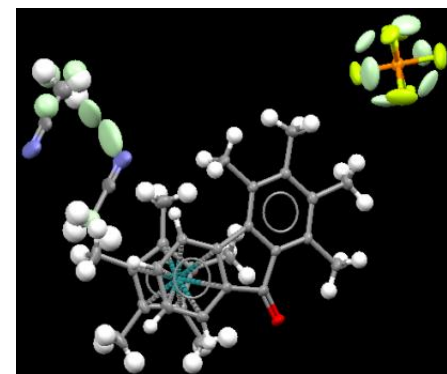
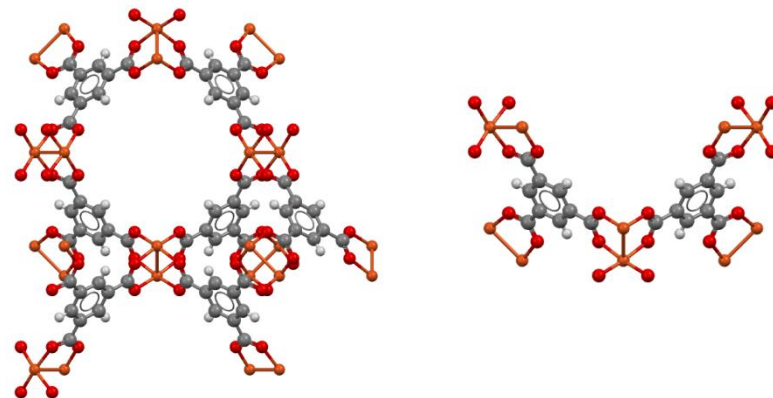
# CSD-Xpedit

- CSD-Xpedit removes internal processing restrictions
  - Large structures
  - Structure Factors
    - Time to mandate?
- Easier access to deposited data
  - Just CCDC number or publication DOI needed
  - Immediate access
- Better structure comparison and duplicate checking
- Persistent identifiers and links from/to other resources
  - Automatic DOI population
- Faster ‘publication’
  - 10-25 days, down to <10 minutes
- Will allow community curation



## Scientific benefits

- Better polymer representation
- ADPs and occupancy factors more visible
- Consistent interpretations





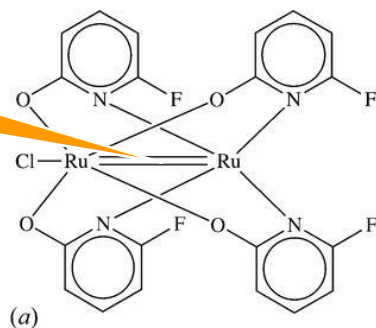
# Curating The Cambridge Structural Database

- Curation
  - Coordinates to 'chemistry'
  - Duplication
  - Families
  - Enrichment
- Distribution and access
  - Searching
  - Analysis
  - Application

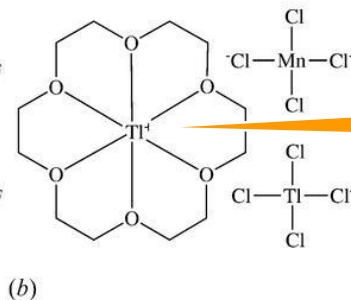


# Chemical structure from atomic coordinates

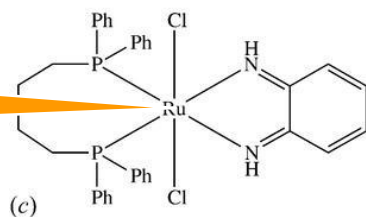
bond order?



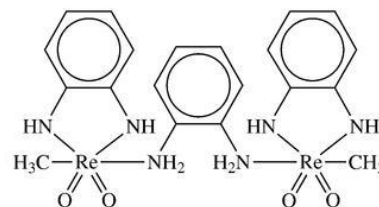
charges?



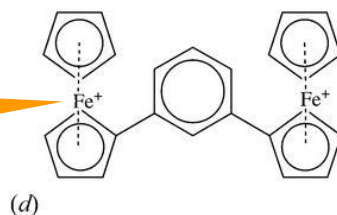
oxidation state?



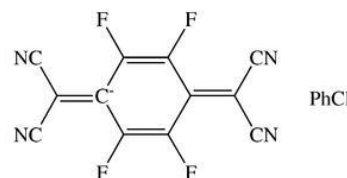
bond type?



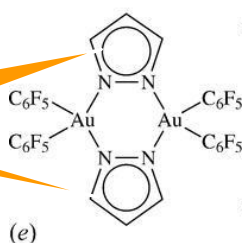
charges?



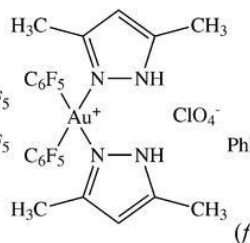
charges?



aromatic?

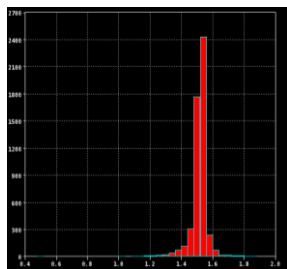


carbene coordinated?

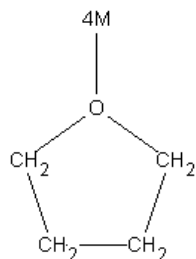




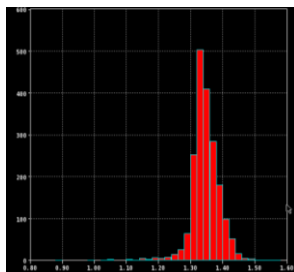
# DeCIFer: Automatic Assignment of Chemistry



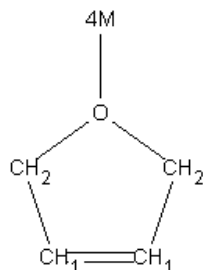
C-C, prob = 0.004



7208 hits



C=C, prob = 0.20



2 hits

Conflicting Evidence

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}.$$

Bayes' Theorem

Low probability bond lengths:

C5-C6 1.405, av(CSD) = 1.505, prob = 0.001

C2-C3 1.345, av(CSD) = 1.514, prob = 0.001

C3-C4 1.338, av(CSD) = 1.514, prob = 0.001

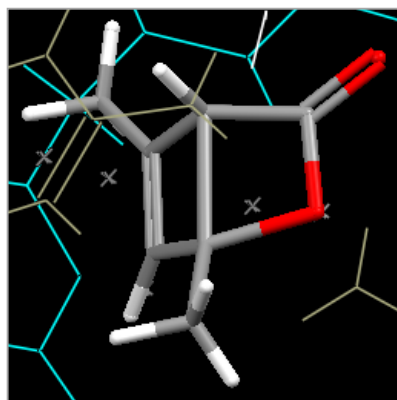
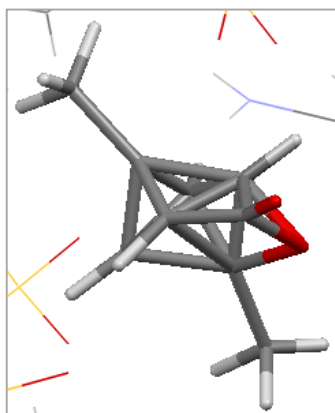
C3-C6 1.798, av(CSD) = 1.546, prob = 0.001

**Reliability level: 2**

Chemical Assignment +  
Reliability Report

*DeCIFer also attempts to automate resolving of disorder and generating diagrams and names*

## Avoiding duplication of effort



Under UV radiation the clathrated pyrone molecule converts to a disordered mixture of square-planar 1, 3-dimethylcyclobutadiene and rectangular-bent 1, 3-dimethylcyclobutadiene in van der Waals contact with a carbon dioxide molecule. The ratio of the square-planar to rectangular-bent 1, 3-dimethylcyclobutadiene clathrate is modelled with occupancies 0.6292:0.3708.

Unresolved disorder

Resolved disorder with editorial comment



## Current challenges

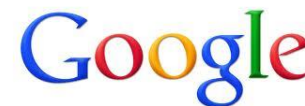
- Increase 'discoverability' of crystal structure data
  - More links – databases, aggregators and publications
  - Improved access to deposited data
- More depositions
- Enriched depositions
  - Diagrams, structure factors
- Community curation





## Improve access to deposited data sets

- The CSD and CSD System are ubiquitous in structural chemistry
- Anyone can access individual deposited data sets
- Access to these data sets must be embedded in other resources
  - Priority to remove technical barriers
  - Establish conditions of use consistent with modern age and desires of rights holders
    - By building on existing services that make data freely available to the community





# Adoption outside small molecule communities

CANCER

## Drug for an 'undruggable' protein

Scientists have long aimed to develop drugs against the cancer-associated protein KRAS, but without success. An approach that targets the oncoprotein's cellular localization reignites lost enthusiasm. [SEE LETTER P.638](#)

NICOLE M. BAKER & CHANNING J. DER

**H**uman *RAS* genes have two claims to notoriety. First, they make up the most frequently mutated oncogene family in human cancer, having a prevalence of one in every three cases<sup>1</sup>. Second, despite more than three decades of intensive effort, no effective pharmacological inhibitor of the *RAS* oncoprotein has reached the clinic. So it is exciting that, on page 638 of this issue, Zimmermann *et al.*<sup>2</sup> report\* the identification and characterization of a small-molecule inhibitor that interferes with the localization of KRAS — the *RAS* isoform most commonly mutated in human cancers — to the plasma membrane surrounding cells<sup>3</sup>.

Following their synthesis in the cytoplasm,

\*This article and the paper under discussion<sup>2</sup> were published online on 22 May 2013.

*RAS* proteins are initially inactive<sup>4</sup>. They then undergo a series of rapid post-translational modifications that ensure their association with the inner leaflet of the plasma membrane, where these proteins exert their normal, as well as their cancer-associated, signalling activity. Therefore, most efforts aimed at anti-*RAS* drug discovery have involved indirect approaches to block the activities of proteins that either promote plasma-membrane association of *RAS* or are components of its downstream signalling pathway.

The key post-translational modification of *RAS* involves the addition of a 15-carbon farnesyl lipid tail in a reaction catalysed by the farnesyltransferase enzyme. This modification facilitates *RAS* association with membranes and is essential for proper *RAS* localization and activity, having prompted intensive efforts in the 1990s to develop

farnesyltransferase inhibitors (FTIs).

Despite promising results in preclinical studies, however, the results of clinical trials with FTIs were disappointing. The inhibitors blocked membrane association of the HRAS isoform, but lacked antitumour activity in cancers involving mutated KRAS (and NRAS). KRAS could still associate with the plasma membrane through an unexpected compensatory activity of the farnesyltransferase-related enzyme geranylgeranyltransferase-1, which modifies *RAS* with a geranylgeranyl, rather than a farnesyl, group. This discouraging outcome greatly dampened interest in targeting *RAS* — and, in particular, its membrane association — for cancer treatment. Instead, ongoing efforts have mainly focused on inhibitors of the RAF-MEK-ERK and the PI3K-AKT signalling cascades downstream of *RAS*.

Zimmermann *et al.* describe an approach aimed at disrupting KRAS membrane association that warrants reassessment of the current strategies. The authors identify and characterize a small-molecule inhibitor of PDEδ, a protein that can bind to and regulate the trafficking of *RAS* and *RAS*-related proteins to membrane compartments<sup>5–8</sup> (Box 1). Specifically, PDEδ contains a deep, hydrophobic pocket capable of binding the lipid moiety of farnesylated proteins, in particular *RAS*.

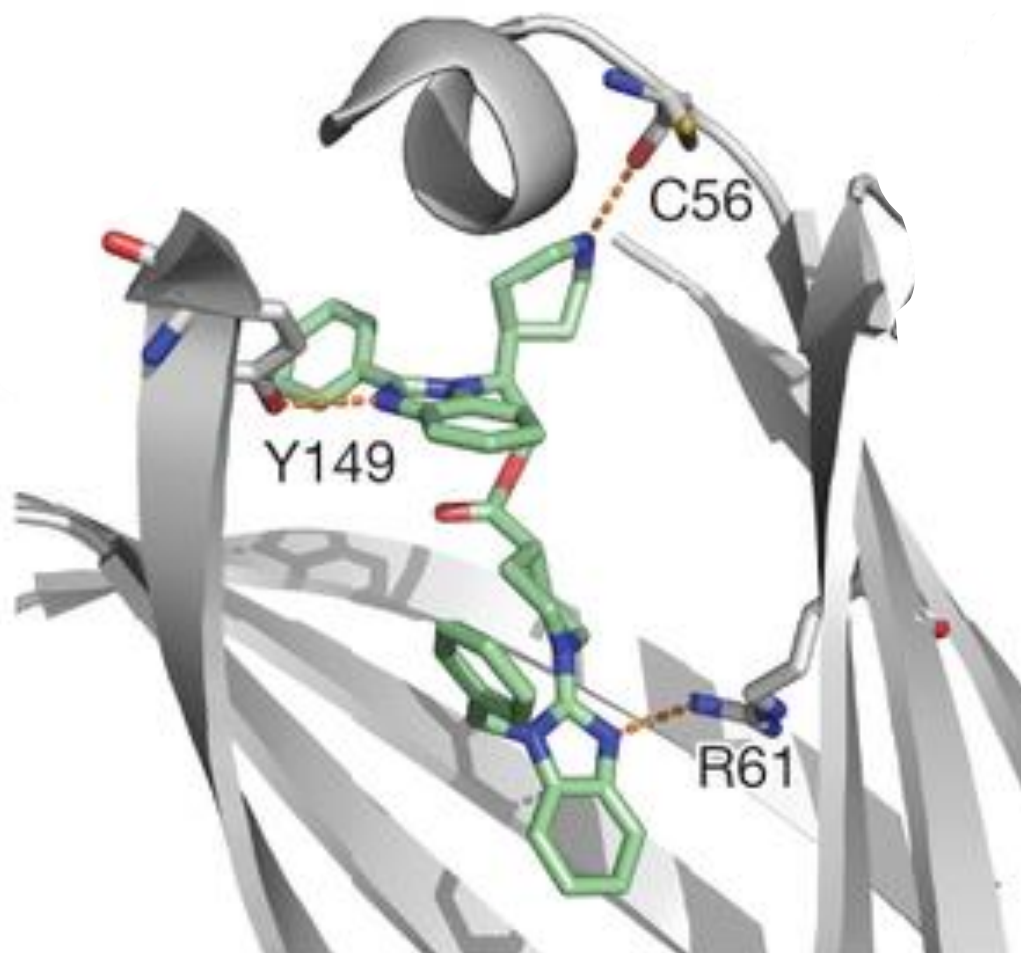
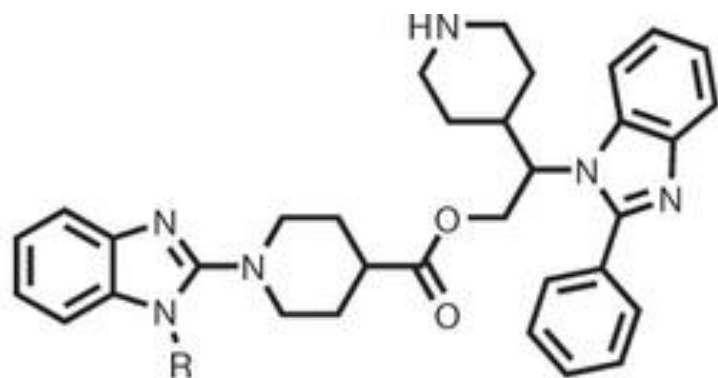
An earlier study<sup>9</sup> found that suppression of PDEδ levels disrupts *RAS* association with the plasma membrane and impairs the growth of *RAS*-mutant cancer cells. This finding

30 MAY 2013 | VOL 497 | NATURE | 577

© 2013 Macmillan Publishers Limited. All rights reserved

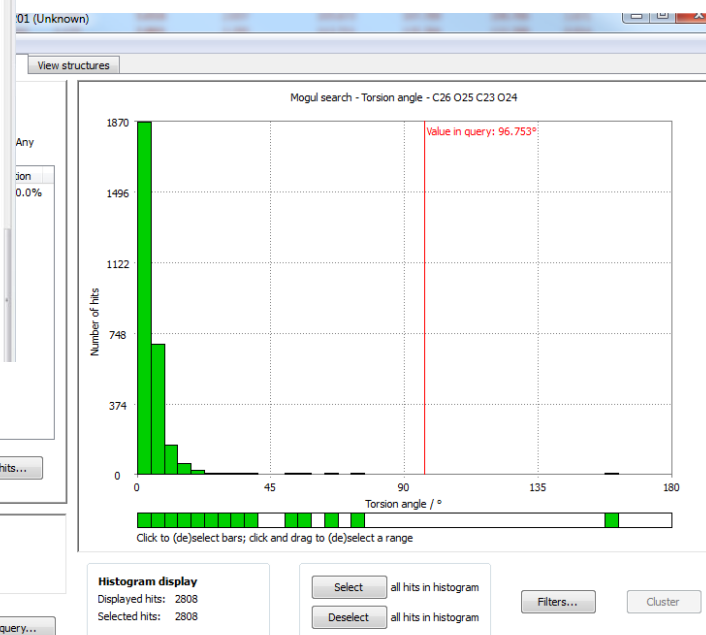
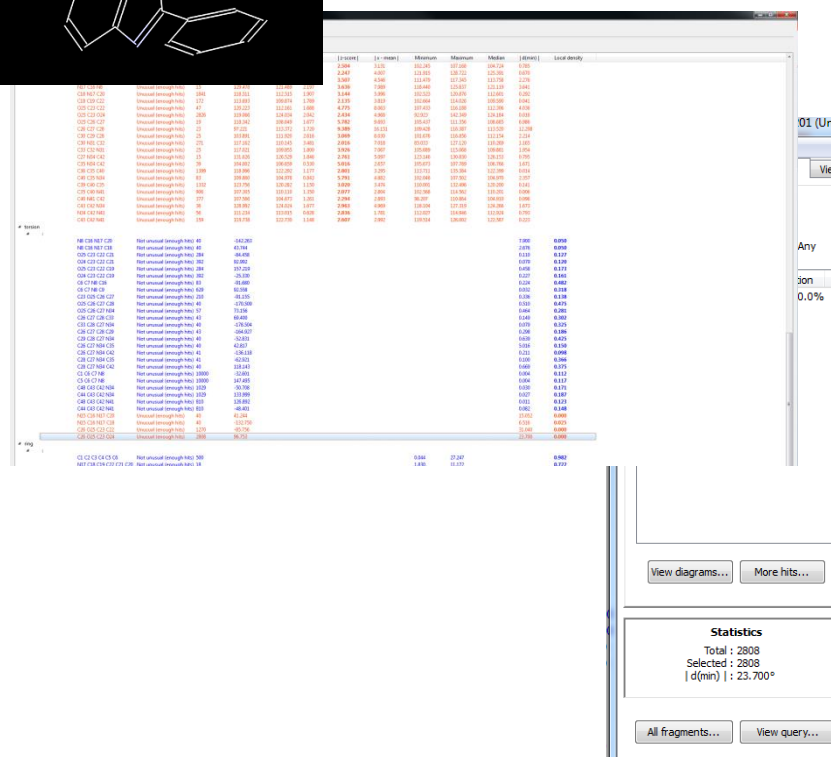
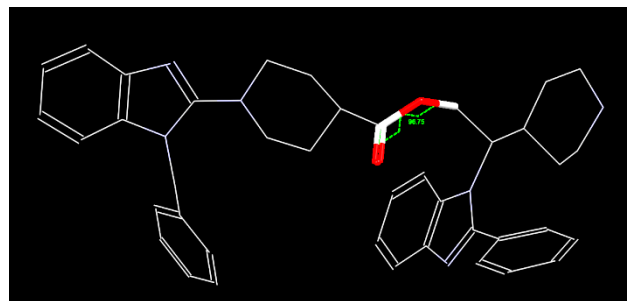


## Adoption outside small molecule communities





# Adoption outside small molecule communities





## Adoption outside small molecule communities

- Structural biology community still struggling with chemistry
- A failure by the small molecule community, including the CCDC to communicate value
  - Some success through collaboration
    - PDB
    - Global Phasing
    - COOT
    - CCP4
- Do current business models contribute to this ‘bad science’?
  - The need for academic contributions restricts sharing
  - Lack of clarity of ‘permissions’



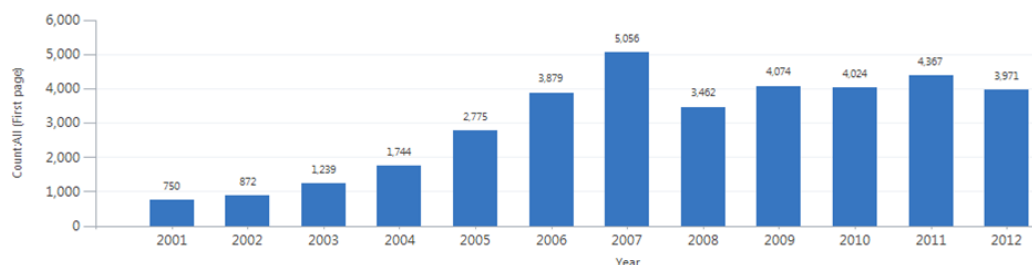
## Current challenges

- Increase ‘discoverability’ of crystal structure data
  - More links – databases, aggregators and publications
  - Improved access to deposited data
- More depositions
- Enriched depositions
  - Diagrams, structure factors
- Community curation



# Structure publishing

- Pre 2000 CCDC and IUCr recognised many valuable structures were never published
  - IUCr launched *Acta Crystallographica Section E* Structure Reports Online
  - 350 (extra?) structures per month



- Recent reduction in structures from *Acta E*
  - No sign yet that they are going elsewhere
- CCDC will assign DOIs to entries
- CSD to be listed in Thomson Reuters Data Citation Index
- What else should we do to encourage 'publication'



## Current challenges

- Increase ‘discoverability’ of crystal structure data
  - More links – databases, aggregators and publications
  - Improved access to deposited data
- More depositions
- Enriched depositions
  - Diagrams, structure factors
- Community curation





## Difficulties with non-semantic data

- CSD Editors still extract data from PDFs
  - Some automatic, some manual
- Utterly stupid way of working
  - Semantic data -> non-semantic, obfuscated data -> semantic data
- Heroic efforts to extract data from PDFs
  - Is this the right place to expend effort?
  - All data could be semantic, especially as supplementary



## Current challenges

- Increase ‘discoverability’ of crystal structure data
  - More links – databases, aggregators and publications
  - Improved access to deposited data
- More depositions
- Enriched depositions
  - Diagrams, structure factors
- Community curation



# Unanswered Questions

- Time to mandate processed data deposition?
  - Who decides?
- Is free access to individual CIFs enough?
- Should the community continue to fund the CSD at point of use
  - Consumer pays
  - This requires access and sharing restrictions
- Should research councils 'pay'
  - Currently the model for many countries
    - Eases move to Open Data, raises questions of sustainability
  - Desperately confused situation in the UK
- Could the entire crystallography community adopt a single *modus operandi*
  - Only Open Access publishing
  - All (processed) data available (semantic form)
  - APCs for papers
  - APCs for data



## Acknowledgements

- Ian Bruno – Strategic Relationships Manager
- Suzanna Ward - CSD Group Manager
- Matthew Lightfoot – Editor in Chief of the CSD
- CCDC staff
- The 253,000 authors of structure containing publications