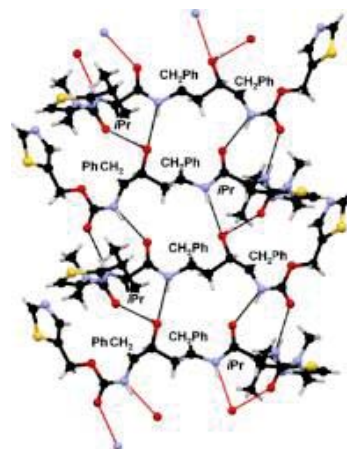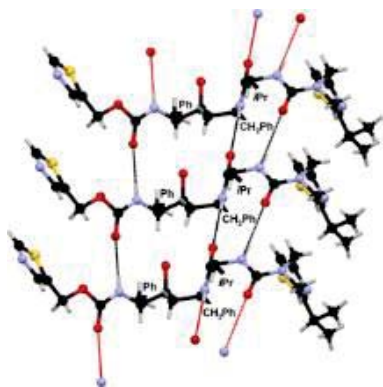# Database-driven discovery

**Suzanna Ward and Eric Rogers**
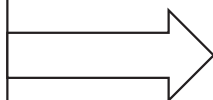
**The Cambridge Crystallographic Data Centre**

# Creation of the CSD

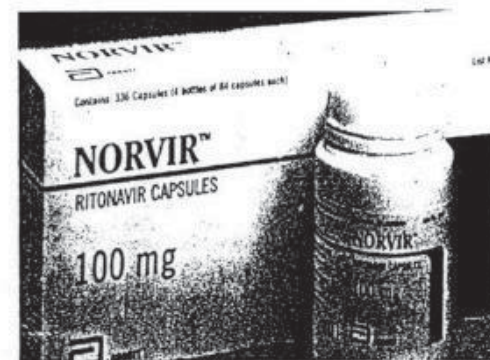# Can Structural Knowledge Mitigate Risk?



Different interactions → Different solubility, Different stability

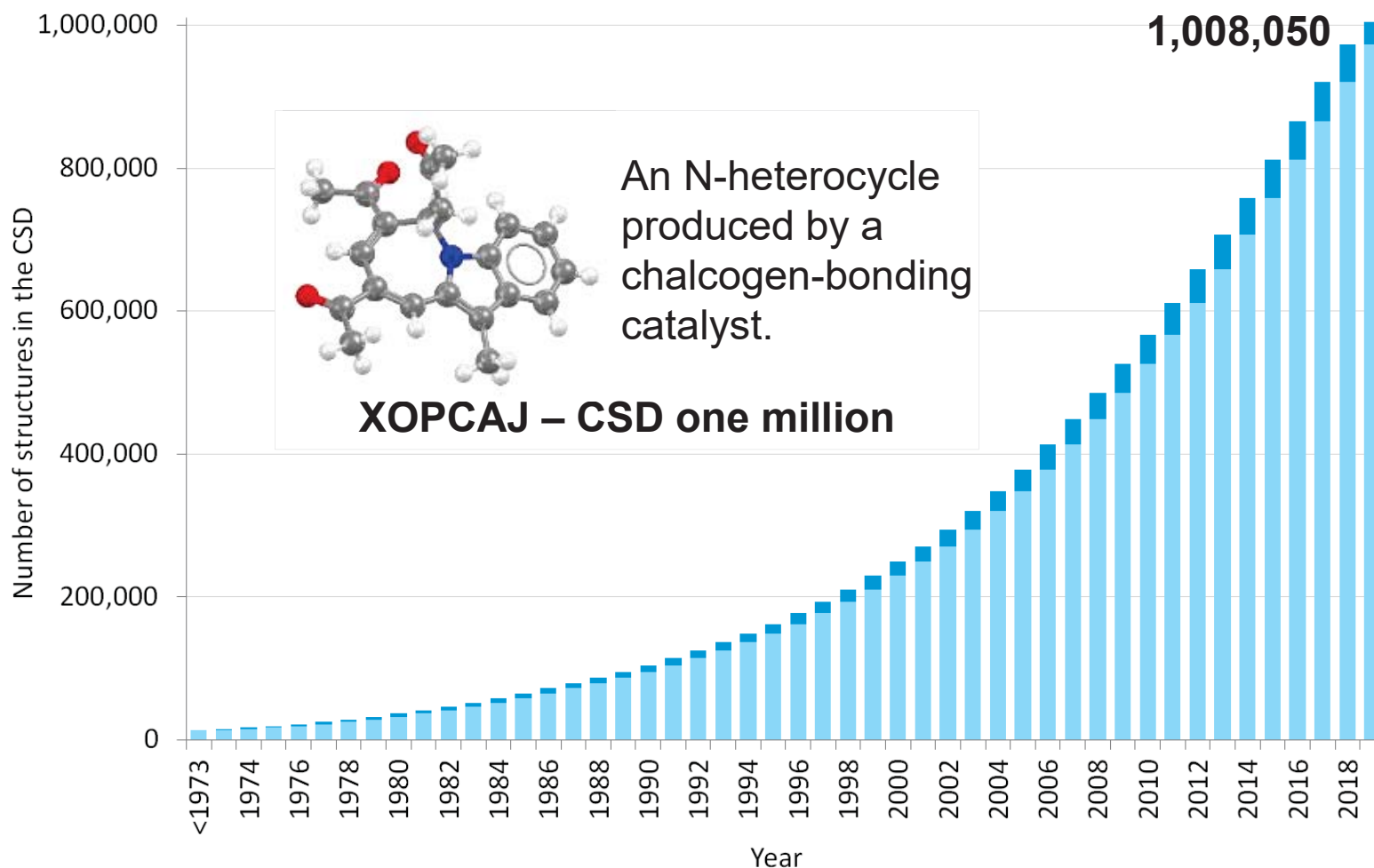Manufacturing problems hit Abbott's HIV drug ritonavir

Capsules of Abbott Laboratories' protease inhibitor Norvir (ritonavir) are likely to become unavailable by the middle of August. The company has a problem with the manufacture of the anti-HIV capsules which it cannot resolve at present.

The problem relates to "undesirable" crystal formation. Abbott says that a series ... ples from a number of marketed batches of capsules were examined and there was no ...

Capsules unlikely to be available from mid-August

# The Cambridge Structural Database (CSD)



An N-heterocycle produced by a chalcogen-bonding catalyst.

**XOPCAJ – CSD one million**

1,008,050

- Every published structure
  - Inc. ASAP & early view
  - *CSD Communications*
  - Patents
  - University repositories

- Every entry enriched and annotated by experts

- Discoverability of data and knowledge

- Sustainable for over 54 years

# CSD One Million

# The Vision

BERNAL'S VISION: FROM DATA TO INSIGHT

by Dr Olga Kennard OBE FRS

THE J D BERNAL LECTURE 1995
delivered at
BIRKBECK COLLEGE, LONDON

We clearly recognised even in those early days, that data banks have three principal functions. Firstly they must gather together existing knowledge and make it readily available to the scientific community. Secondly they can be used to reduce a large number of observations to a small set of constants and rules, and in this way transform a data base to a knowledge base. Such a knowledge base may obviate the need for further individual experiments in specific areas. Thirdly, they facilitate the comparison and collective analysis of individual results to gain insight into new or as yet unexplained phenomena. These ideas have been at the heart of the work of the Cambridge Crystallographic Data Centre and the driving force for improving methods of data collection, storage and dissemination. Most importantly they influenced development of computer programs and methodologies which are needed for the analysis and transformation of the accumulated information. (5)

# Inside the CSD

**Organic 43%**

**Metal-Organic 57%**

At least one transition metal, lanthanide, actinide or any of Al, Ga, In, Tl, Ge, Sn, Pb, Sb, Bi, Po

**Not Polymeric 89%**

**Polymeric: 11%**

**Single Component 56%**

**Multi Component 44%**

- **Organic**
  - Drugs
  - Agrochemicals
  - Pigments
  - Explosives
  - Protein ligands

- **Metal-Organic**
  - Metal Organic Frameworks
  - Models for new catalysts
  - Porous frameworks for gas storage
  - Fundamental chemical bonding

- **Additional data**
  - 10,860 polymorph families
  - 169,218 melting points
  - 840,667 crystal colours
  - 700,002 crystal shapes
  - 23,622 bioactivity details
  - 9,740 natural source data
  - > 250,000 oxidation states

- **Links/subsets**
  - Drugbank
  - Druglike
  - MOFs
  - PDB ligands
  - PubChem
  - ChemSpider
  - Pesticides

**Why Is Our AI Revolution Built On Free Data Rather Than Good Data?**

**Kalev Leetaru** Contributor ⓘ
AI & Big Data
*I write about the broad intersection of data and society.*

Getty Images.   GETTY

One of the greatest challenges confronting the modern AI revolution

Many of the most pressing challenges facing AI today revolve around its poor-quality training data. Bias, brittleness, ease of fooling, lack of representational edge case examples to fall back upon: all of these key problems trace their roots at least in part to poor quality training data. While algorithmic improvements could help, so too could having proper training data.

# AI and Machine Learning

- AI and machine learning techniques are evolving rapidly

- But the consequences of using poor quality data can be far reaching
  - Incorrect scientific conclusions
  - Wasted investment and effort
  - A loss of trust
  - Ultimately poor business decisions.

https://www.forbes.com/sites/kalevleetaru/2019/02/05/why-is-our-ai-revolution-built-on-free-data-rather-than-good-data/

# Curating the CSD

- Each dataset expertly curated

- Datasets enhanced
  - Chemical connectivity
  - Compound names
  - 2D chemical diagrams
  - Additional experimental data
  - Bibliographic information

# Depositing the Data

# Guidelines

## The CCDC CIF Deposition Guidelines

When preparing your CIF for deposition please include as much information as possible and check it carefully. This is especially true for *CSD Communications* where there is no paper to describe the chemistry and experimental details leading to your structure. If you choose to publish your data as a *CSD Communication* please remember to provide all the authors/crystallographers/chemists who contributed to the crystallographic experiment as authors of the data. If we are unable to validate your structure from the information you have provided we may contact you. If we cannot resolve the issue, unfortunately, we may not be able to add your structure to the CSD.

🇨🇳 Guidelines in Chinese

All experimental CIF files (including those from powder diffraction experiments) should contain an R-factor. This should be consistent with the crystallography being performed correctly and to the best ability that would be expected from the material and equipment used. We would like all experimental CIFs to contain:

- R-factors (R1, wR2, Rint)
- GooF
- Shift/ESD (to show that the refinement has converged)
- Explanation of any problems with numbers of reflections and parameters
- Any residual electron density
- Details of squeeze/solvent masking
- Atomic Displacement Parameter (ADP) values
- Temperature – cell and data collection temperatures match
- Experimental set up including mounting device and instrument type
- HKL included
- RES included

We would encourage you to take advantage of the IUCr checkCIF reports built in to the deposition page. This can highlight issues to check with your structure that can be clarified in the validation reply form, particularly in the case of A- or B- level alerts. Ideally, treatment of disorder or partial occupancy atoms should be clear and of course, no non-positive definite atoms!

To allow us to create the most accurate representation of your structure please provide as much additional information on the "Enhanced Data" page as is appropriate for your structure. Some chemical issues we commonly encounter when processing data into the CSD are:

- Given formula and crystal formula don't agree. Particular attention should be paid to hydrogen atoms which may not be located in the experiment. It would be very helpful to us to have a complete moiety formula (including unlocated hydrogens and any SQUEEZE/MASK species not located, if known)
- Charge balance, particularly for variable metal oxidation states and radicals
- Missing hydrogen atoms, especially on oxygen atoms that could be hydroxy/oxy/aqua ligands and for polyoxometalate structures
- Unusual bonding, tautomers or metal-metal bonding
- Poorly handled or unmodelled disorder
- Unexplained void space not accounted for by SQUEEZE or MASK procedures

Further information that will benefit the users of your structure and that will enable the correct identification of any previous versions of your structure are:

- Stereochemical determination method, if relevant
- Crystallisation solvent/conditions
- Melting point
- Details of re-refinement – please tell us if the structure is a re-refined version of an existing CSD entry.
- Refcodes or CCDC numbers of any known related structures; i.e. by temperature / stereochemistry/ pressure; e.g. "high temperature determination of REFCODE"

If you have any further queries, please contact us via our Enquiries Page.

# Adoption by the Community



% CIFs containing HKL data in CSD deposits

# What Else Could We Do?

- **Improved peer review**
  - Mandate crystallographic review of all structure-containing papers
  - Educate reviewers on nature of CheckCIF alerts

- **File requirements**
  - CIF + structure factors
  - Refinement instructions?
  - CheckCIF report?

- **Validation checks**
  - CheckCIF integration
  - Unit cell checks (with HKL checks? Or chemistry check?)
  - Geometry analysis?

- **Additional files available to reviewers?**

# Joint CSD and ICSD Services



Over 180,000 entries from the Inorganic Crystal Structure Database (ICSD) now available through Access Structures

Joint Access

Joint Deposition

# Curation and Chemistry Assignment

# Using the CSD to Help With Curation

An automated probabilistic approach using data in the CSD

```
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_adp_type
_atom_site_occupancy
_atom_site_symmetry_multiplicity
_atom_site_calc_flag
_atom_site_refinement_flags
_atom_site_disorder_assembly
_atom_site_disorder_group
Cl1 Cl 0.5993(2) 1.0007(7) 0.8131(17) 0.044(3) Uani 0.50 1 d PDU A 1
S1 S 0.5321(3) 0.8260(6) 0.9322(3) 0.0327(11) Uani 0.50 1 d PDU A 1
C2 C 0.5529(4) 0.8802(9) 0.8184(9) 0.029(4) Uani 0.50 1 d PDU A 1
C3 C 0.5286(7) 0.8174(18) 0.7440(7) 0.031(4) Uani 0.50 1 d PDU A 1
H3A H 0.5350 0.8343 0.6771 0.037 Uiso 0.50 1 calc PR A 1
C4 C 0.4918(8) 0.7220(19) 0.7783(8) 0.027(4) Uani 0.50 1 d PDU A 1
C5 C 0.4900(6) 0.7171(14) 0.8779(9) 0.029(4) Uani 0.50 1 d PDU A 1
C12 Cl 0.3202(2) 0.4982(6) 1.0830(5) 0.0586(15) Uani 0.50 1 d PDU A 1
S2 S 0.38755(19) 0.6658(5) 0.9578(5) 0.0400(10) Uani 0.50 1 d PDU A 1
```

**Assignment of chemistry is required to make data findable, interoperable and reusable**

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

# Challenges



Missing atoms



Element assignment



Disorder



Poor geometries

# The Human Touch

- Each entry looked at by expert Scientific Editors

- Reliability scores focusses editorial efforts

- Manual validation of automated chemical interpretations improves automated methods

https://www.ccdc.cam.ac.uk/Community/blog/CSD-data-curation-the-human-touch/

# Revisiting Data

Targeted improvements allow improved integrity, consistency, discoverability and value of data


Ensure standardisation of early CSD entries


Creation and maintenance of subsets


Enrichment of data

Oxidation states

# Melting Points in the CSD

>170,000 Melting Points

Study Temperature relative to MP
*before* additional CCDC validation

Study Temperature relative to MP
*after* additional CCDC validation

# Maintaining Data Integrity

- **Integrity** – Completeness, consistency and trustworthiness

- **Data completeness** – Trends in reporting of metadata
  - Identify CSD Deposit checks and enhancements
  - Identify new filters to allow CSD users to better select fit for purpose data

- **Consistency** – Looking at experimental metadata to identify trends in information supplied

- **Trustworthiness** – Establishing automatic identification of potential cases of misconduct – including fraudulent and plagiarised data

# Following Standard Ethical Practises

- CCDC is now a Member of the **C**ommittee **o**n **P**ublication **E**thics.

- COPE's objective is "*to educate and advance knowledge in methods of safeguarding the integrity of the scholarly record for the benefit of the public*".

- Membership gives us access to COPE resources and COPE advice – helping us deal with publication ethics and data integrity and issues.



C|O|P|E

Member since 2019
AC00039

https://publicationethics.org/about/governance

# Making Crystallographic Data FAIR

## CCDC database workflows

# FAIR Data Principles

# FAIR Data Principles

**F**indable 🔍
**A**ccessible
**I**nteroperable ⚙
**R**eusable ♻

"all research objects should be Findable, Accessible, Interoperable and Reusable (FAIR) both for machines and for people" (Wilkinson, M. D. et al., 2016 : 3)



https://www.force11.org/group/fairgroup/fairprinciples



http://www.datafairport.org/

# FAIR Data Policies and Guidelines

## Plan S

"Although the Plan S principles refer to peer-reviewed scholarly publications, cOAlition S also strongly encourages that research data and other research outputs are made as open as possible and as closed as necessary."

https://www.coalition-s.org/principles-and-implementation/

## European Commission

"..the implementation of FAIR data needs to go hand-in-hand with the principle that data created by publicly-funded research must be as Open as possible and as closed as necessary. The EC and Member States should consider FAIR and Open as complementary concepts and address both in policy. " *(EU Commission, 2018 : 10)*

# Funder Research Data Sharing Policies

"Since 2017, all Horizon 2020 projects are part of the Open Research Data Pilot by default. The Principal Investigator must:

• Develop a data management plan in the first 6 months of the project and keep it up-to-date throughout their project;
• Deposit their research data in a suitable research data repository;
• Make sure third parties can freely access, mine, exploit, reproduce and disseminate their data;
• Make clear what tools will be needed to use the raw data to validate research results, or provide the tools themselves."

"ERC beneficiaries are encouraged to take part in the H2020 Open Research Data Pilot, but this is not compulsory.
Those who take part in the Open Research Data Pilot must adhere to the obligations outlined above."

https://www.data.cam.ac.uk/funders

# Publisher research data policies

"The Royal Society of Chemistry believes that, where possible, all data associated with the research in a manuscript should be **freely available** in an **accessible** and **usable** format, enabling other researchers to replicate and build on that research. Therefore, in addition to providing the data required for submission (as detailed above) we encourage authors to deposit as much data as possible that is related to the research in their article. This should be in appropriate and **publicly available repositories**"

https://www.rsc.org/journals-books-databases/journal-authors-reviewers/prepare-your-article/experimental-data/

"It is the practice of IUCr journals to provide **free access** to all supplementary materials and supporting data files deposited with a published article."

https://journals.iucr.org/services/authorrights.html

# Aspects of FAIR Data

## Findable

- Globally unique and persistent identifiers
- Rich metadata descriptions
- (Meta)data available in a searchable resource

## Interoperable

- Standard formats for representation
- Use of FAIR vocabularies
- References to other (meta)data

## Accessible

- (Meta)data retrievable by their identifier
- Standard, open communication protocols
- Metadata accessible even when data are not

## Reusable

- Described with a plurality of attributes
  - data usage licenses
  - detailed provenance
  - domain-relevant community standards

# Data Repositories

Role of Repositories in making data FAIR:

- ❑ Long-term data preservation
- ❑ Continued access to data
- ❑ Assignment of identifiers and DOIs
- ❑ Added descriptive metadata
- ❑ Searchable databases
- ❑ Data curation

Repository frameworks and systems

### DOI assignment



### Open Archival Information System (OAIS) Reference Model



### Global and domain specific membership groups



### Data preservation policies & plans

### Trusted Repository Certification



Digital Preservation Coalition and Brian Lavoie (2014) The Open Archival Information System (OAIS) Reference Model: Introductory Guide (2nd Edition), Digital Preservation Coalition. DOI: 10.7207/twr14-02

# Data Repositories

How to find a repository for your research data

❑ Search a repository registry



https://www.re3data.org/

❑ Use a repository recommended by the publisher or funding body

❑ Search for repositories among accreditation bodies



https://www.coretrustseal.org/



https://www.icsu-wds.org/services/certification

# Crystallographic Information File: CIF

International Union of Crystallography

Commission on Crystallographic Data
Commission on Journals
Working Party on Crystallographic Information

**The Crystallographic Information File (CIF): a New Standard Archive File for Crystallography\***

By Sydney R. Hall

*Crystallography Centre, University of Western Australia, Nedlands 6009, Australia*

Frank H. Allen

*Crystallographic Data Centre, University Chemical Laboratory, Lensfield Road, Cambridge CB2 1EW, England*

and I. David Brown

*Institute for Materials Research, McMaster University, Hamilton, Ontario L8S 4M1, Canada*

A standard format for archive and exchange of crystallographic data
- derived model
- processed data (structure factors)
- metadata about raw data (imgCIF)

A standard format for archive and exchange of crystallographic data
- derived model
- processed data (structure factors)
- metadata about raw data (imgCIF)

```
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_adp_type
_atom_site_occupancy
_atom_site_symmetry_multiplicity
_atom_site_calc_flag
_atom_site_refinement_flags
_atom_site_disorder_assembly
_atom_site_disorder_group
C11 Cl 0.5993(2) 1.0007(7) 0.8131(17) 0.044(3) Uani 0.50 1 d PDU A 1
S1 S 0.5321(3) 0.8260(6) 0.9322(3) 0.0327(11) Uani 0.50 1 d PDU A 1
C2 C 0.5529(4) 0.8802(9) 0.8184(9) 0.029(4) Uani 0.50 1 d PDU A 1
C3 C 0.5286(7) 0.8174(18) 0.7440(7) 0.031(4) Uani 0.50 1 d PDU A 1
H3A H 0.5350 0.8343 0.6771 0.037 Uiso 0.50 1 calc PR A 1
C4 C 0.4918(8) 0.7220(19) 0.7783(8) 0.027(4) Uani 0.50 1 d PDU A 1
C5 C 0.4900(6) 0.7171(14) 0.8779(9) 0.029(4) Uani 0.50 1 d PDU A 1
C12 Cl 0.3202(2) 0.4982(6) 1.0830(5) 0.0586(15) Uani 0.50 1 d PDU A 1
S2 S 0.38755(19) 0.6658(5) 0.9578(5) 0.0400(10) Uani 0.50 1 d PDU A 1
```

# CIF as a FAIR data format

**Findable**

- Searchable fields for identifiers and metadata descriptions

**Interoperable**

- Standard dictionary and vocabularies
- Standard format for processed and derived data

**Accessible**

- Trusted searchable data repositories:
- Cambridge Structural Database
- Inorganic Crystal Structure Database
- Protein Data Bank

**Reusable**

- Data provenance
- Software packages and parameters
- Quality metrics

# CCDC Dataset Workflow

Deposition (Pre-ingest) → Deposition (Ingest) → Publication → Curation → Data Access

## Data Deposition (Pre-ingest)

Data deposited by data producer

CIF, HKL and FCF data deposited via the CCDC deposition and validation service

**CCDC checks run on deposited files:**

- Structure factor check
- IUCr checkCIF
- Unit Cell Check

Metadata and terms of deposition confirmed by the depositor

Depositor's responses to checks and checkCIF reports added to deposit record.

Data submitted

Option for depositor to retrieve deposited files

## Data Deposition (Ingest)

CCDC automatic validation and duplicate check of data files

Does data pass automatic validation?

No → Data manually checked by CCDC staff

Does data pass manual validation?

Yes → Data assigned Deposition Number(s)

No → Email sent to depositor requesting additional information/data

Yes → Data stored for long-term preservation

Email sent to depositor containing Deposition Number(s)

## Data Publication

Pre-publication metadata communicated to CCDC by journal publishers

Full publication metadata communicated to CCDC by journal publishers

Publication information updated via journal scanning by CCDC staff

Pre-publication metadata added to data record

Full publication metadata added to data record

## Data Access

Data made accessible to publishers and referees pre-publication (dependent on checks)

Deposited data made freely accessible from Access Structures

Data accessible pre-publication to depositors via My Structures

Data enters the CSD

## Data Curation

Data transferred to FIZ for curation into the ICSD

Does structure meet criteria for curation into the CSD?

No

Yes → Structure enters queue for scientific validation

Structure validated by CCDC's scientific editors

Remarks added to entry and Refcode assigned / confirmed

Accessible for download at: https://www.ccdc.cam.ac.uk/Community/depositastructure/scientific-data-preservation/

# Data Deposition (Pre-ingest)

**Data Deposition (Pre-ingest)**

Data deposited by data producer

CIF, HKL and FCF data deposited via the CCDC deposition and validation service

**CCDC checks run on deposited files:**

- Structure factor check
- IUCr checkCIF
- Unit Cell Check

Metadata and terms of deposition confirmed by the depositor

Depositor's responses to checks and checkCIF reports added to deposit record.

Option for depositor to retrieve deposited files

Data submitted

- Manual User Actions:
  - ❑ Provide personal details and upload files
  - ❑ Fix syntax errors (if any)
  - ❑ Explain why no structure factor data
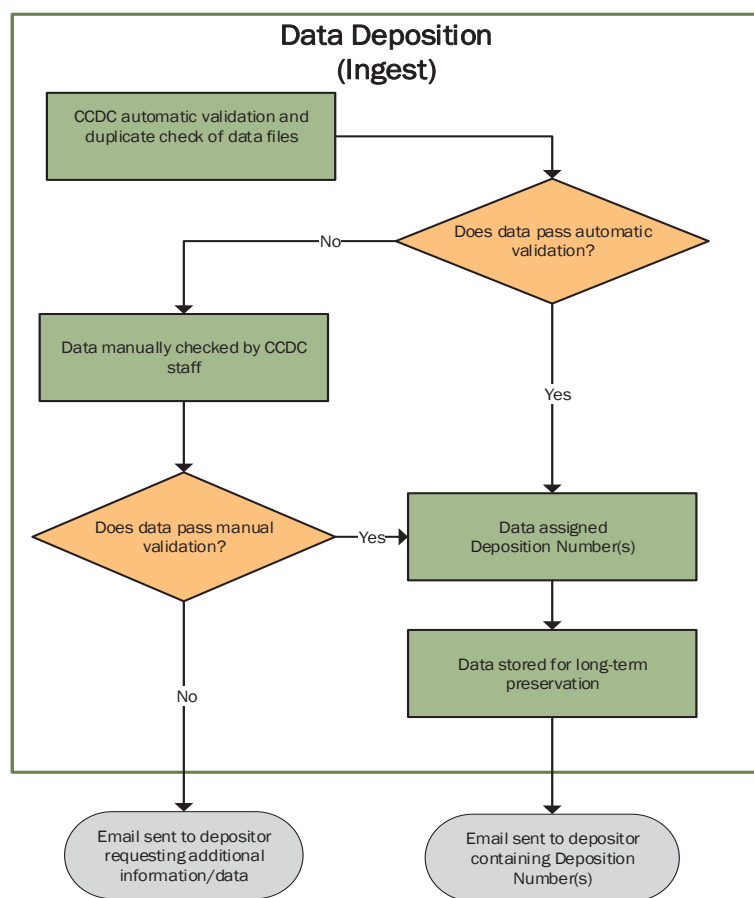  - ❑ Add explanations for checkCIF alerts
  - ❑ Add crystallographer details
  - ❑ Provide known publication details
  - ❑ Add additional scientific metadata
  - ❑ Review and confirm

- Automated Actions:
  - ❑ Syntax check
  - ❑ Generation of checkCIF report
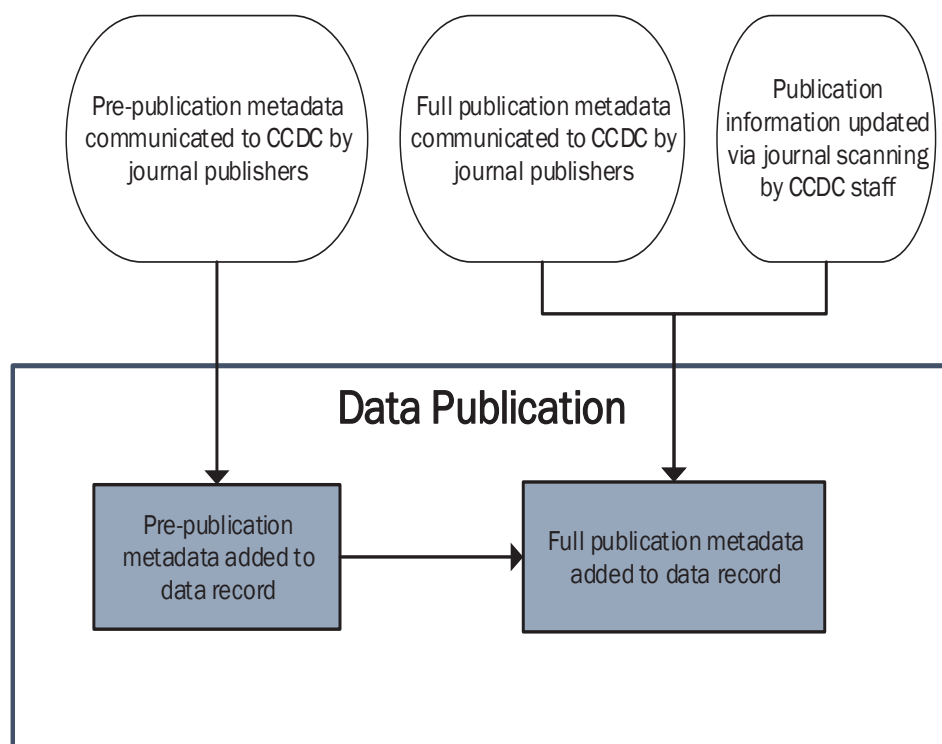
# Data Deposition (Ingest)



**Data Deposition (Ingest)**

- CCDC automatic validation and duplicate check of data files
- Does data pass automatic validation?
  - No → Data manually checked by CCDC staff
    - Does data pass manual validation?
      - No → Email sent to depositor requesting additional information/data
      - Yes → Data assigned Deposition Number(s)
  - Yes → Data assigned Deposition Number(s)
- Data assigned Deposition Number(s) → Data stored for long-term preservation → Email sent to depositor containing Deposition Number(s)

- Automated Actions:
  - ❑ Syntax check
  - ❑ Check for duplicate dataset
  - ❑ Internal record creation
  - ❑ Assigning identifiers

- **Manual CCDC Actions:**
  - ❑ Dealing with non-standard file formats
  - ❑ Investigating duplicate datasets
  - ❑ Updating records with resubmitted data

# Data Publication



Pre-publication metadata communicated to CCDC by journal publishers

Full publication metadata communicated to CCDC by journal publishers

Publication information updated via journal scanning by CCDC staff

## Data Publication

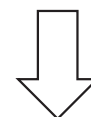Pre-publication metadata added to data record

Full publication metadata added to data record

## Tabernabovines A−C: Three Monoterpenoid Indole Alkaloids from the Leaves of *Tabernaemontana bovina*

Yang Yu,[†,‡] Mei-Fen Bao,[†,‡] Jing Wu,[†,‡] Jing Chen,[†,‡] Yu-Rong Yang,[†,§] Johann Schinnerl,[‖] and Xiang-Hai Cai[*,†,§]

### Accession Codes

CCDC 1916676 contains the supplementary crystallographic data for this paper. These data can be obtained free of charge via www.ccdc.cam.ac.uk/data_request/cif, by emailing data_request@ccdc.cam.ac.uk, or by contacting The Cambridge Crystallographic Data Centre, 12 Union Road, Cambridge CB2 1EZ, UK; fax: +44 1223 336033.
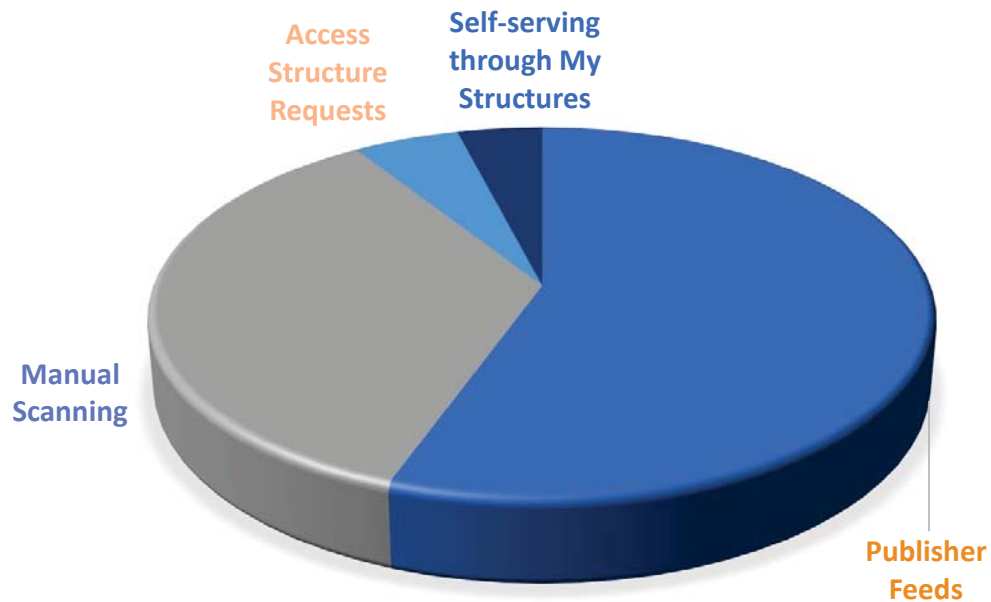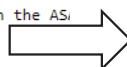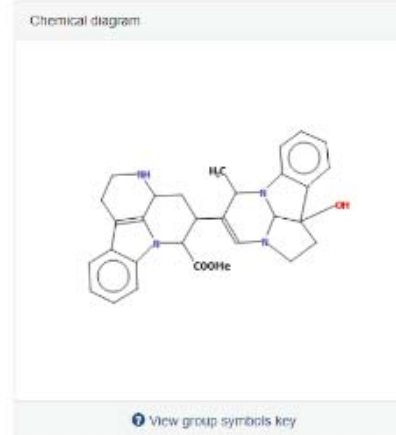


CCDC

# Data Publication

**PUBLICATION INFO SOURCES, ESTIMATE, MARCH 2019**



Pie chart segments labeled: Access Structure Requests, Self-serving through My Structures, Publisher Feeds, Manual Scanning

**Sources of Publication Information:**

❑ Pre-publication metadata communicated by journal publisher feeds

❑ Full publication metadata communicated and updated by journal publisher feeds

**Manual CCDC Actions:**

❑ Reviewing publication details

❑ Publication information updated via journal scanning by CCDC staff

❑ Publication information communicated by researchers wanting to access data

# Data Publication



**Accession Codes**

CCDC **1916676** contains the supplementary crystallographic data for this paper. These data can be obtained free of charge via www.ccdc.cam.ac.uk/data_request/cif, by emailing data_request@ccdc.cam.ac.uk, or by contacting The Cambridge Crystallographic Data Centre, 12 Union Road, Cambridge CB2 1EZ, UK; fax: +44 1223 336033.

```xml
<?xml version="1.0" encoding="UTF-8"?>
<rss xmlns:acs="http://distiller.acs.org/iapps/wld/rss/ccdc/module" version="2.0">
  <channel>
    <title>CCDC Feed</title>
    <link>https://pubs.acs.org</link>
    <description>This feed contains all of the CIF codes for all articles that are in the AS
    <language>en-US</language>
    <copyright>Copyright 2019 American Chemical Society</copyright>
    <lastBuildDate>Tue, 16 Jul 2019 10:18:15 GMT</lastBuildDate>
    <item>
      <title>Tabernabovines A-C: Three Monoterpenoid Indole Alkaloids from the Leaves of &lt;i
      <link>http://dx.doi.org/10.1021/acs.orglett.9b02060</link>
      <pubDate>Thu, 11 Jul 2019 04:00:00 GMT</pubDate>
      <author>Yang Yu†‡, Mei-Fen Bao†‡, Jing Wu†‡, Jing Chen†‡, Yu-Rong Yang†§, Johann Schinne
      <guid isPermaLink="false">10.1021/acs.orglett.9b02060</guid>
      <acs:issn-print>1523-7060</acs:issn-print>
      <acs:issn-electronic>1523-7052</acs:issn-electronic>
      <acs:ccdc>1916676</acs:ccdc>
    </item>
```

# Data Access

## Data Access

- Data made accessible to publishers and referees pre-publication (dependent on checks)

- Deposited data made freely accessible from Access Structures

- Data accessible pre-publication to depositors via My Structures

- Data enters the CSD

- Data is made accessible:
  - ❑ Pre-publication to reviewers and depositors to facilitate preparation of manuscripts
  - ❑ Immediately through Access Structures once data is published
  - ❑ Through the CSD once curated by CCDC's scientific editors.

- Processes for making data more findable and interoperable:
  - ❑ Identifiers added to data entries
  - ❑ Links from publication articles to data created
  - ❑ Links from CCDC datasets to other databases

# Standard Identifiers and Interoperability

Data should be considered legitimate, citable products of research…

https://www.force11.org/datacitation

Dataset Publication
CCDC 610092: Experimental Crystal Structure Determination. **A. Crystallographer**, *Cambridge Crystallographic Data Centre* (2007)
http://dx.doi.org/10.5517/ccngvdb

- The CCDC registers DOIs for datasets through DataCite
- Metadata for CCDC datasets is openly accessible via DataCite
- Foundation for interoperability and formalising data citation

10.5517/CCPHZ37

ORCID IDs for Researchers

At least 30% of current CSD depositors provide an ORCID ID

Andrew Bond
ORCID ID
https://orcid.org/0000-0002-1744-0489

# Links from Articles to CCDC Data

# Making data accessible though the CSD

Data not published in a scientific journal can be curated into the CSD and made available to the community as a **CSD Communication**

Structures from your PhD **thesis** can be made publicly available through the CSD.
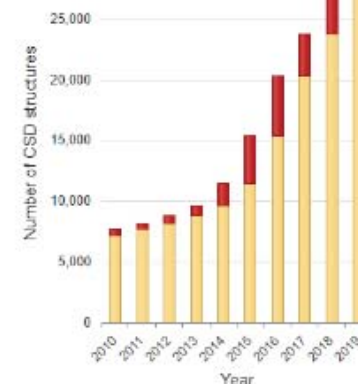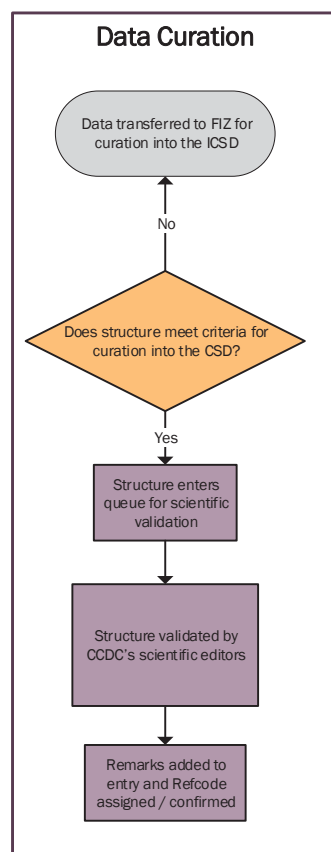


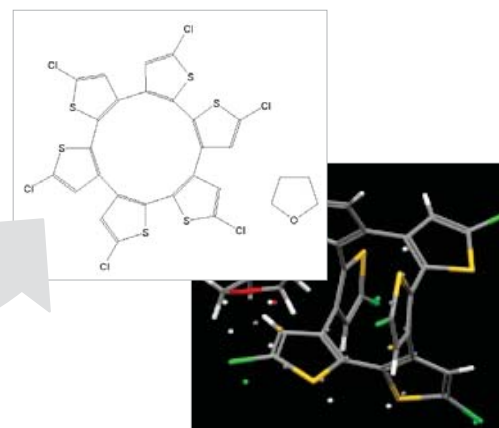https://www.ccdc.cam.ac.uk/Community/csd-communications/

# Data Curation

```
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_adp_type
_atom_site_occupancy
_atom_site_symmetry_multiplicity
_atom_site_calc_flag
_atom_site_refinement_flags
_atom_site_disorder_assembly
_atom_site_disorder_group
C11 C1 0.5993(2) 1.0007(7) 0.8131(17) 0.044(3) Uani 0.50 1 d PDU A 1
S1 S 0.5321(3) 0.8260(6) 0.9322(3) 0.0327(11) Uani 0.50 1 d PDU A 1
C2 C 0.5529(4) 0.8802(9) 0.8184(9) 0.029(4) Uani 0.50 1 d PDU A 1
C3 C 0.5286(7) 0.8174(18) 0.7440(7) 0.031(4) Uani 0.50 1 d PDU A 1
H3A H 0.5350 0.8343 0.6771 0.037 Uiso 0.50 1 calc PR A 1
C4 C 0.4918(8) 0.7220(19) 0.7783(8) 0.027(4) Uani 0.50 1 d PDU A 1
C5 C 0.4900(6) 0.7171(14) 0.8779(9) 0.029(4) Uani 0.50 1 d PDU A 1
C12 C1 0.3202(2) 0.4982(6) 1.0830(5) 0.0586(15) Uani 0.50 1 d PDU A 1
S2 S 0.38755(19) 0.6658(5) 0.9578(5) 0.0400(10) Uani 0.50 1 d PDU A 1
```

## Data Curation

- Data transferred to FIZ for curation into the ICSD
  - No
- Does structure meet criteria for curation into the CSD?
  - Yes
- Structure enters queue for scientific validation
- Structure validated by CCDC's scientific editors
- Remarks added to entry and Refcode assigned / confirmed

Assignment of chemistry is required to make data findable, interoperable and reusable

- A reliable chemical representation is essential for enabling reuse and application of crystallographic data

- Representation is generated at CCDC using a combination of automated processes and manual validation

# Enablers of FAIR Crystallographic Data

- **Standard formats and identifiers**
  - ❏ Crystallographic Information Format and dictionaries (CIF)
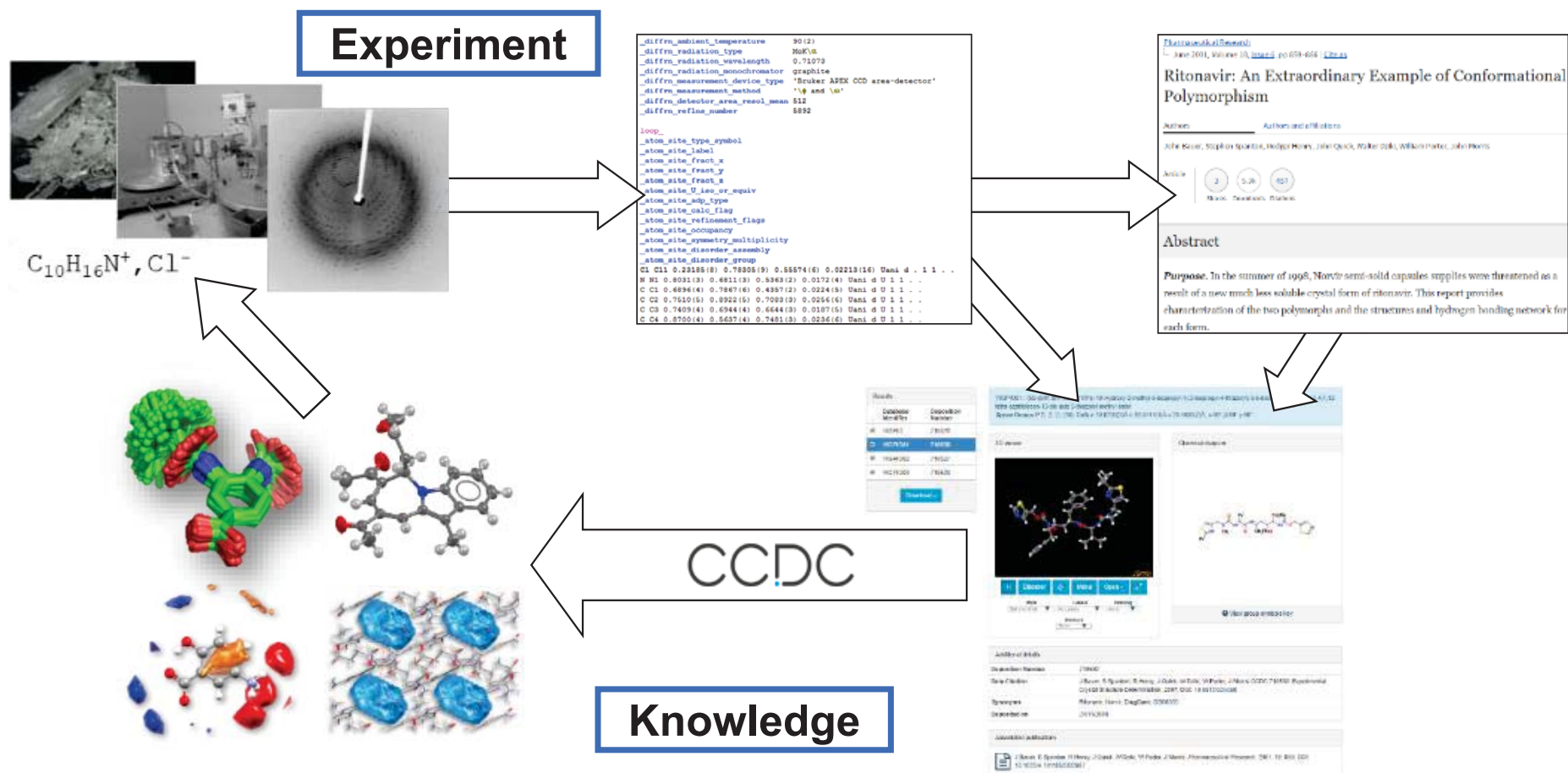  - ❏ Standard Identifiers and associated infrastructure (DOIs, ORCID, InChI…)
- **Community stakeholders**
  - ❏ Instrument providers and software developers adopting standards
  - ❏ Publishers and editors encouraging use of standards for publication
  - ❏ Repositories and databases providing access to enriched data
  - ❏ International Unions supporting and promoting standards
  - ❏ Individual researchers and others championing research data standards
- **Tools and services that make it easy to make data FAIR**

# From Data to Knowledge
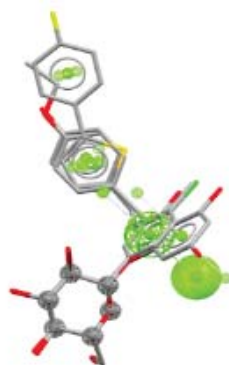
# From Experiment to Knowledge
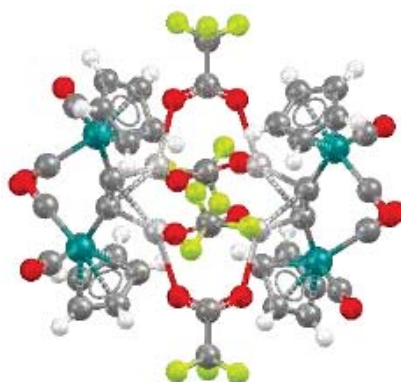
# Advanced Services and Software



**CSD-Enterprise**
All CCDC applications and software
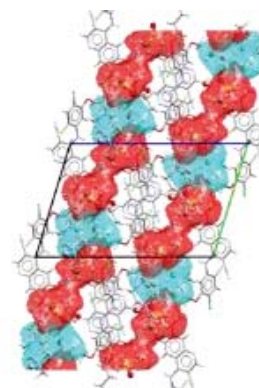
**CSD-Discovery**

*To discover new molecules with pharmaceutical applications*

**CSD-System**

*To search, visualise, analyse and communicate structural data*

**CSD-Materials**

*To understand and predict solid form stability and properties*

**The Cambridge Structural Database**

# Generating Insights

- The **CSD Python API** enables you to create tailored scripts using full array of CSD functionality

- Answer targeted research questions or integrate access with other software

- Functions include:
    - Full search capabilities
    - Geometry analysis
    - Interaction analysis
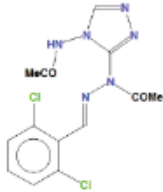    - Descriptor calculation
    - 2D diagram generation

# Increasing Complexity

- Increasing:

  - Formula weights
  - Unit Cells
  - Number of elements

# Trends in Experimentation

# Elements in the CSD



Percentage of structures that contain each element from before (red) and after (yellow) 2009

# Drugs

- Top 200 Pharmaceutical Products
  - **By retail sales in 2018**
  - Produced by the Njarðarson Group
  - The University of Arizona
  - Drugs already in the CSD coloured green
- *J. Chem. Ed.* **2010**, *87*, 1348

# Identifying Trends in Drug Structures



https://www.ccdc.cam.ac.uk/Community/blog/insights-into-drug-like-compounds-from-crystal-data/

# The CSD and the PDB

## Linking

- Between CSD and PDB ligands



icode: TOBBOB

a: $C_3H_2F_6O$

Name: 2-Difluoromethoxy-1,1,1,2-tetrafluoroethane

onym: (+)-Desflurane; PDB Chemical Component code: DSF; DrugBank: DB01189

CCDC Class:

Source:

Melting Point: 147K

Colour: colorless

Extra Information: absolute configuration; inhalation anaesthetic activity

## CSD-CrossMiner

- Pharmacophore query tool
- Searches the CSD and PDB

# Using the Data



CSD FINWEE10/PDB FK5

Match PDB ligands to best representative CSD molecules

# Solid Form Informatics

- The term "**solid form informatics**" first introduced in mid-2000s by Bob Docherty (Pfizer):

  *Use of structural knowledge to inform key decisions in pharmaceutical development*

- Solid form informatics now a key part of the solid form development workflow at most major pharmaceutical companies



**Molecule** → **Form** → **Particle**

# From Data to Knowledge

Individual data points from different datasets combine to provide information that aids in the discovery and optimisation of new chemical entities



Molecular Shape

Molecular Interactions

Taylor *et al. J. Chem. Inf. Model.*, (2014) 54 (9), 2500. Wood, P. A. *et al. CrystEngComm* (2013) **15**, 65

# Predicting Unlikely Interactions

Predictive analytics is used to identify the likelihood of specific molecular interactions occurring from similar crystal structures



Galek *et al, CrystEngComm,* (2009), 11, 2634 - 2639

The integration of solid-form informatics into solid-form selection

Neil Feeder[a], Elna Pidcock[a], Anthony M. Reilly[a], Ghazala Sadiq[a], Cheryl L. Doherty[b], Kevin R. Back[b], Paul Meenan[c] and Robert Docherty[b]

One in half a million: a solid form informatics study of a pharmaceutical crystal structure

Peter T. A. Galek,*[a] Elna Pidcock,[a] Peter A. Wood,[a] Ian J. Bruno[a] and Colin R. Groom[a]

Navigating the Solid Form Landscape with Structural Informatics

Peter T. A. Galek, Elna Pidcock, Peter A. Wood, Neil Feeder, Frank H. Allen

Book Editor(s): Yuriy A. Abramov

Knowledge-based H-bond prediction to aid experimental polymorph screening

Peter T. A. Galek,*[ab] Frank H. Allen,[a] László Fábián[ab] and Neil Feeder[c]

# CSD-Materials: Targeted Solutions



Crystal Packing Similarity

Motif Search &
Packing Feature Search

Conformer Generator

Full Interaction Maps

Complex Structural Analysis

Solid Form Risk Assessment

Hydrogen Bond Propensity

DASH

Solid Form Design

Hydrate Analyser &
Solvate Analyser

Calculations

Molecular Complementarity

# Individual value


16d

"A **search** of the Cambridge Structural **Database** using a series of **pharmacophore queries** led to the discovery of an O-spiroketal C-arylglucoside scaffold. Subsequent chemical examination combined with computational modelling resulted in the identification of the clinical candidate 16d (CSG452, tofogliflozin), which is currently under phase III clinical trials."

Yoshihito Ohtake et al *Journal of Medicinal Chemistry* 2012 *55* (17), 7828-7840  (Roche, Chugai)

# Collective value

## Articles

# The Supramolecular Synthon Approach to Crystal Structure Prediction

J. A. R. P. Sarma*,† and Gautam R. Desiraju*,‡

gvk bioSciences Pvt. Ltd., #210, 'My Home Tycoon', 6-3-1192, Begumpet, Hyderabad 500 016, India, and School of Chemistry, University of Hyderabad, Hyderabad 500 046, India

W This paper contains enhanced objects available on the Internet at http://pubs.acs.org/crystal.

**ABSTRACT:** A new approach has been proposed for the ab initio crystal structure prediction of small organic molecules. This exercise forms a part of the recent blind test on crystal structure prediction conducted by the Cambridge Crystallographic Data Centre. The method uses as a starting point lists of low energy structures generated by an exhaustive computational procedure, namely, the Polymorph Predictor program in $Cerius^2$. Such computational procedures take into account only the enthalpic factors in crystallization. A further difficulty is that information relating to crystallization kinetics is very hard to obtain directly. However, such kinetic information is implicitly contained in the experimental structures that are found in crystallographic databases. Therefore, in our approach, the low energy structures obtained in the Polymorph Predictor program are reranked after consideration of experimental structures of structurally similar molecules. Operationally, this is most conveniently carried out after identification of possible supramolecular synthons in the Cambridge Structural Database. These synthons are representative structural units that convey critical information that relates isolated molecules with their resulting crystal structures. Of the three molecules in the blind test, the present approach was fully successful for one, but only of limited utility in the two others. Reasons for this variability of success are given.

# Collective value

we note that if the CSD were to be significantly larger than what it is today, say, around a million refcodes, CSP with the synthon-based approach could be successfully employed for a much wider variety of molecules.

The Su
Structu

J. A. R. P. Sarma*,† and Gautam R. Desiraju*,‡

gvk bioSciences Pvt. Ltd. #210, 'My Home Tycoon' 6-3-1192, Begumpet, Hyderabad 500 016, India, a

Received I

Ⓦ This p

ABSTRA
molecule
Cambrid
by an exh
procedur
relating
contained
the low
experime
identificat
representativ

In summary, and from the viewpoint of CSP, the utilization of structural information could provide a more effective sieve toward the correct solution. As the amount of structural information in crystallographic databases increases, structure prediction would gradually move toward fingerprinting.

crystal structures. Of the three molecules in the blind test, the present approach was fully successful for one, but only of limited utility in the two others. Reasons for this variability of success are given.
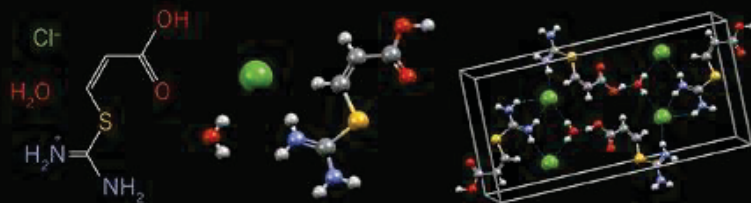
# Using the Collection



CCDC Blind Test Showcases Major Advance in Crystal Structure Prediction Methods

– November 03, 2015

The Cambridge Crystallographic Data Centre (CCDC) announces that the results of its 6th blind test of crystal structure prediction methods demonstrate significant advancement in crystal structure prediction methods in comparison with previous tests.
of polymorphs, salts and hydrate
experimental structures were pre

**CRYSTAL CHALLENGE**
The 3D structure that a molecule adopts in a crystal is very difficult to predict — but defines what properties the molecule has.

| The structural formula of a molecule reveals which atoms are connected at a 2D level. | Chemists are making progress at predicting how complex molecules will assemble in 3D space — there are millions of possibilities. | The 3D orientation repeats in a crystalline lattice with a structure that dictates the molecule's mechanical, chemical and physical properties. |

## nature
International weekly journal of science
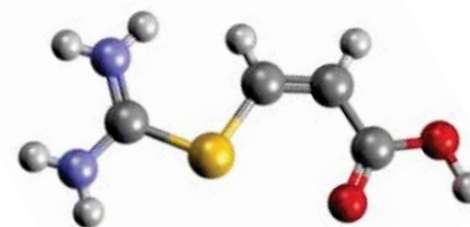
Software predicts slew of fiendish crystal structures

Chemists succeed at forecasting how complex molecules will assemble in 3D.

**Elizabeth Gibney**

vember 2015

n the structure of an organic molecule on a napkin and it may not be apparent that there are s of possible ways that it could assemble as a 3D crystal. Now, a collaboration of dozens of sts and computer programmers has successfully predicted the crystal structure of five, ex, 'drug-like' organic molecules — using nothing but a 2D map showing which atoms ct to which.

chievement, annou      27 Octob      workshop in Cambridge, UK, paves the way for are that would cut      cost of the des      manufacture of drugs and other chemical cts, as well as further our understanding of fundamental chemistry.

# Thank you