



An Introduction to
The Cambridge Crystallographic Data Centre
The Cambridge Structural Database
Chemistry Data Initiatives

Ian Bruno

Cambridge Crystallographic Data Centre

@ijbruno @ccdc_cambridge



The Cambridge Crystallographic Data Centre

International Data Repository

Archive of crystal structure data
High quality scientific database

Scientific Software Provider

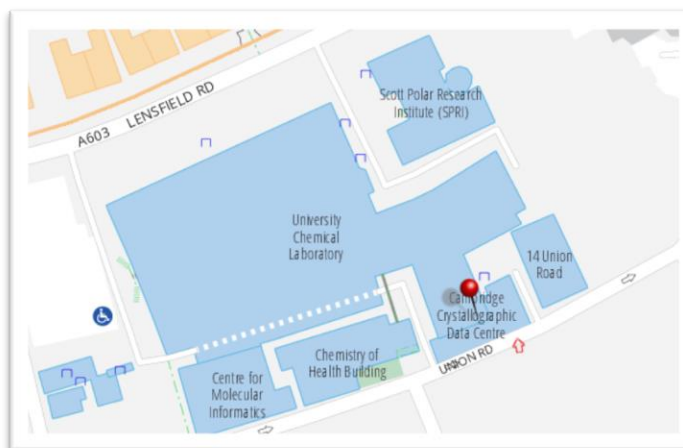
Search/analysis/visualisation tools
Scientific applications

Collaborative Research Organisation

New methodologies
Fundamental research

Education and Outreach

Conferences, Workshops,
Training, Teaching





The Crystallographic Data Centre Cambridge



Olga Kennard, Founding
Executive Director

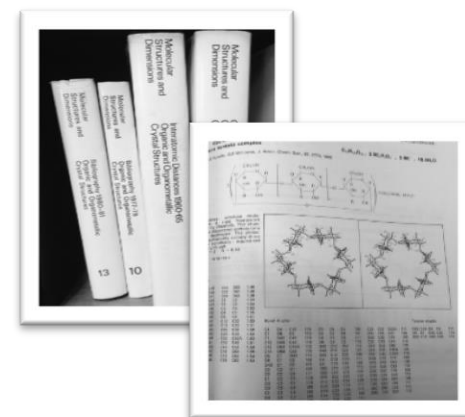
“The database was established in 1965 to fulfil a dream of myself and a great scientist, the polymath J.D. Bernal. We had a passionate belief that the collective use of data would lead to the discovery of new knowledge which transcends the results of individual experiments.”

Kennard, O. “From Private Data to Public Knowledge.” *The Impact of Electronic Publishing on the Academic Community*. Ed. I Butterworth. Portland Press Ltd, 1997. 159–166.

Established 1965 with funding from the Royal Society

Main objective was to assemble a computer-based file of information and data

Data was first published in printed volumes generated from the computer file





Early Publication of Crystal Structures

Hand-typed tables of coordinates in journal articles manually transcribed into database records

2178 J. CHEM. SOC. DALTON TRANS. 1985

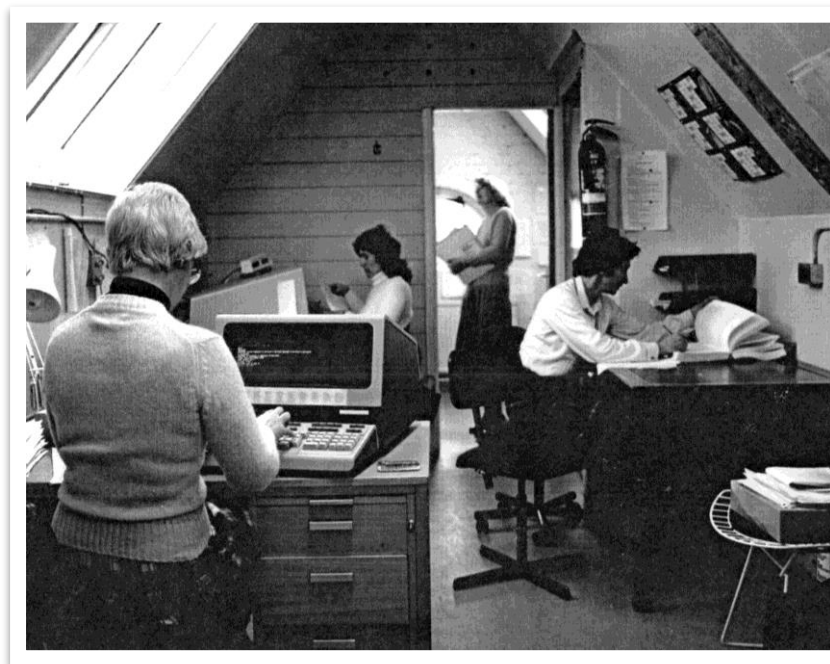
Table 1. Crystallographic data and details of data collection and processing for $ML(NO_2)_2$, with M = Cu [in (1)], Ni [in (2)], and Cd [in (3)]

	(1)	(2)	(3)
Stoichiometry			
M			
Lattice type			
Space group			
a/Å			
b/Å			
c/Å			
α°			
β°			
γ°			
$U/\text{Å}^3$			
Z			
$D_x/\text{g cm}^{-3}$			
F(000)			
$\mu(\text{Mo-K}\alpha)/\text{cm}^{-1}$			
Approximate crystal dimensions (mm)			
Number of settings			
θ range/ $^\circ$ (cell dimensions)			
h range			
k range			
l range			
Number of reflexions measured			
independent			
observed			
Final R			
Final R			

J. CHEM. SOC. DALTON TRANS. 1985 2179

Table 2. Atomic co-ordinates with estimated standard deviations in parentheses

Atom	X/a	Y/b	Z/c	Atom	X/a	Y/b	Z/c
(a) Compound (1) ($\times 10^3$ for Cu, $\times 10^4$ for others)							
Cu	7 734(2)	11 519(5)	62 931(2)	C(27)	2 251(3)	671(8)	4 844(4)
N(01)	1 493(1)	-1 223(3)	6 454(2)	C(30)	2 044(2)	483(5)	6 905(2)
C(10)	1 359(2)	-1 289(5)	6 758(3)	N(31)	1 888(1)	1 149(3)	7 431(2)
N(11)	708(1)	-1 425(3)	6 566(2)	N(32)	1 386(1)	1 906(3)	7 246(2)
N(12)	379(1)	-356(3)	6 561(2)	C(33)	1 395(2)	2 512(4)	7 809(2)
C(13)	-162(2)	-764(4)	6 588(2)	C(34)	1 907(2)	2 157(5)	8 333(2)
C(14)	-171(2)	-2 078(4)	6 570(2)	C(35)	2 220(2)	1 302(4)	8 094(2)
C(15)	386(2)	-2 475(4)	6 586(2)	C(36)	918(3)	3 421(6)	7 819(3)
C(16)	-651(2)	142(7)	6 546(4)	C(37)	2 788(2)	622(6)	8 413(3)
C(17)	650(3)	-3 767(5)	6 606(4)	N(40)	-304(2)	2 265(5)	5 590(3)
C(20)	1 549(2)	-339(5)	5 777(2)	O(41)	136(1)	2 399(3)	6 163(2)
N(21)	1 509(2)	902(3)	5 488(2)	O(42)	-717(2)	3 031(5)	5 486(3)
N(22)	1 099(1)	1 738(3)	5 555(2)	O(43)	-272(2)	1 385(5)	5 217(2)
C(23)	1 123(2)	2 745(4)	5 177(2)	N(50)	1 960(2)	3 011(5)	3 346(2)
C(24)	1 539(3)	2 539(6)	4 854(2)	O(51)	1 968(2)	4 072(5)	3 604(3)
C(25)	1 786(2)	1 388(5)	5 035(2)	O(52)	1 526(2)	2 301(2)	3 201(2)
C(26)	743(3)	3 859(6)	5 139(3)	O(53)	2 417(2)	2 613(5)	3 337(4)
(b) Compound (2) ($\times 10^4$)							
Ni	5 110(1)	3 357(1)	7 658(1)	N(31)	2 270(6)	1 531(4)	7 437(3)
N(01)	4 516(6)	1 860(4)	6 658(3)	N(32)	2 763(6)	2 672(4)	7 868(3)
C(10)	5 885(9)	995(6)	7 036(5)	C(33)	1 964(8)	2 706(6)	8 528(4)
N(11)	6 172(6)	936(4)	8 049(4)	C(34)	897(8)	1 577(7)	5 516(5)
N(12)	6 420(6)	2 047(5)	8 549(3)	C(35)	1 109(8)	827(6)	7 812(4)
C(13)	6 802(8)	1 740(7)	9 452(5)	C(36)	2 223(9)	3 805(7)	9 165(5)
C(14)	6 818(11)	478(8)	9 522(6)	C(37)	3261(10)	-423(6)	7 459(6)
C(15)	6 405(9)	-15(6)	8 625(5)	N(40)	7 521(8)	5 107(5)	8 260(4)
C(16)	7 092(10)	2 741(8)	10 200(5)	O(41)	6 084(6)	4 914(4)	8 518(3)
C(17)	6 208(11)	-1 311(6)	8 281(6)	O(42)	6 610(7)	5 982(5)	8 553(4)
C(20)	4 684(10)	2 395(6)	5 790(5)	C(20)	7 666(5)	4 302(4)	7 690(3)
N(21)	3 777(7)	3 513(5)	5 656(4)	N(50A)	1 902(9)	8 557(6)	5 048(5)
N(22)	3 928(6)	4 195(5)	6 453(4)	O(51A)	3 077(7)	7 942(6)	4 339(5)
C(23)	2 971(9)	6 137(6)	6 121(6)	O(52A)	3 133(10)	8 780(8)	5 628(7)
C(24)	2 311(10)	5 059(8)	5 134(6)	O(53A)	573(14)	8 512(6)	5 431(6)
C(25)	2 805(9)	3 985(7)	4 857(5)	N(50B)	2 032(9)	8 882(6)	4 874(5)
C(26)	2 759(11)	6 066(7)	6 797(6)	O(51B)	1 353(7)	8 323(6)	4 160(5)
C(27)	2 488(11)	3 386(8)	3 912(5)	O(52B)	3 281(10)	9 440(8)	5 053(7)
C(30)	2 652(8)	1 313(6)	6 551(4)	O(53B)	1 492(14)	8 544(6)	5 645(6)



Kleywegt et al. (1985) *J. Chem. Soc., Dalton Trans.*, 2177-2184
doi:10.1039/DT9850002177



Publication of Crystal Structures Today

Electronic data files deposited and disseminated via the Web and linked with journal articles

```
S11 C13 C14 C15 C16 0.1(4) . . .  
C13 C14 C15 C16 0.1(4) . . .  
C14 C15 C16 C17 0.2(5) . . .  
C13 Cu2 Cu 1.000000 0.500000  
C13 C11 C12 diffn_reflns_point  
C15 _reflns_number_total  
C15 loop _reflns_number_gt  
C13 _atomb _reflns_threshold_ex  
C7 _atomb _reflns_Friedel_cove  
C1 _atomb _reflns_Friedel_frac  
C13 _atomb _reflns_Friedel_frac  
C7 _atomb _reflns_Friedel_frac  
C1 _atomb _reflns_special_deta  
C24 _atomb ;  
S11 Cu1 {Reflections were mer  
C16 S11 {class for the calcul  
C20 N1 0.  
C21 N2 0. _reflns_Friedel_frac  
C20 N4 0. Friedel pairs measur  
S11 C1 0. possible theoretical  
C22 C2 0. systematic absences.  
C3 0. ;  
_sh C4 0. ;  
C5 0. ;  
; _computing_data_coll  
TIT C7 0. _computing_cell_refi  
C8 0. _computing_data_redu  
C9 0. _computing_structure  
CEL C9 0. _computing_structure  
ZER C10 0. _computing_molecul  
LAT C11 0. _computing_publicati  
SYM C12 0. ;  
SFN C19 {loop_  
UNT C20 {atom_site_label  
SIZ C21 {atom_site_type_symb  
TEM C22 {atom_site_fract_x  
L S C25 {atom_site_fract_y  
ACT C26 {atom_site_fract_z  
BOND Cu2 {atom_site_U_iso_or  
CONN C11 {atom_site_adp_type  
LIS _atom_site_occupancy  
FMA _geom _atom_site_site_symb  
PLA ; _atom_site_calc_flag  
WGH All 4 _atom_site_refinemen  
FVA Are 4 _atom_site_refinemen  
COI into _atom_site_refinemen  
and _atom_site_disorder  
CL1 used _atom_site_disorder  
treat Fe1 Fe 0.23371(2) 0.  
CL2 ;  
C11 C1 0.16365(3) 0.  
C12 C1 0.20883(4) 0.  
S11 loop S11 S1 0.38708(3) 0.  
_geom N1 N 0.30132(5) 0.58  
N1 _geom N2 N 0.26923(5) 0.48  
_geom N3 N 0.40083(11) 0.2  
N2 _geom C1 C 0.35204(11) 0.5  
_geom C2 C 0.38286(12) 0.4  
Cu1 C1 H2A H 0.417817 0.404  
Cu1 C3 C 0.36301(14) 0.5065(3) 0.79468(13) 0.0400(7) Uani 1 1 d . . .  
Cu1 H3A H 0.382956 0.476446 0.831310 0.048 Uiso 1 1 calc R U . . .  
Cu1 C4 C 0.31380(14) 0.5843(3) 0.78383(13) 0.0405(7) Uani 1 1 d . . .
```

Issue 20, 2018



From the journal:
Dalton Transactions

The coordination chemistry of the neutral pyridyl silicon ligand [PhSi(6-Me-2-py)₃]

Alex J. Plajer,^a Annie L. Colebatch,^a Markus Enders,^b Álvaro García-Romero,^c &

Dominic S. Wright^{*a}

Author affiliations

Abstract

Difficulties in the preparation of neutral ligands of the type [RSi(2-py)₃] (pyridyl ring unit) have thwarted efforts to expand the coordination chemistry simply switching the pyridyl substituents to 6-methyl-pyridyl groups (allowed smooth, high-yielding access to the [PhSi(6-Me-2-py)₃] ligand coordination chemistry with transition metals. The synthesis, single-crystal dynamics of the new complexes [(PhSi(6-Me-2-py)₃]Cu(CH₃CN)[PF₆], [(PhSi(6-Me-2-py)₃]FeCl₂, [(PhSi(6-Me-2-py)₃]Mo(CO)₃ and [(PhSi(6-paramagnetic Fe²⁺ and Co²⁺ complexes show strongly shifted NMR resonances due to large Fermi-contact shifts. However, magnetic anisotropy also contributes to shifts so that both contributions have to be included in the paramagnetic

Previous Article

Next Article

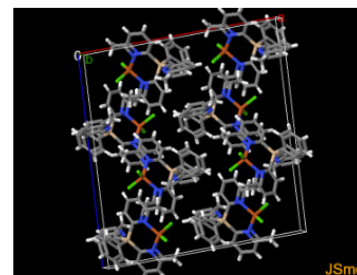
Results

Database Identifier	Deposition Number
<input checked="" type="checkbox"/> LEVYET	1833560
<input checked="" type="checkbox"/> LEVYIX	1833561
<input checked="" type="checkbox"/> LEVYOD	1833562
<input checked="" type="checkbox"/> LEVYUJ	1833563
<input checked="" type="checkbox"/> TIGWUE	1833558
<input checked="" type="checkbox"/> TIGXAL	1833559

Download

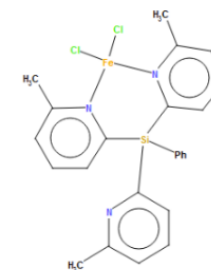
LEVYET : dichloro-[2,2'-(6-methylpyridin-2-yl)(phenyl)silanedyl]bis(6-methylpyridine))-iron
Space Group: C 2/c (15), Cell: a 20.9935(5)Å b 10.2893(2)Å c 22.4986(5)Å, α 90° β 93.9540(10)° γ 90°

3D viewer



H Disorder Menu Open JSmol
Style Labels Packing Measure
Capped Sticks No Labels Unit Cell None

Chemical diagram



View group symbols key

Additional details

Deposition Number	1833560
Data Citation	Alex J. Plajer, Annie L. Colebatch, Markus Enders, Álvaro García-Romero, Andrew D. Bond, Dominic S. Wright CCDC 1833560: Experimental Crystal Structure Determination, 2018, DOI: 10.5517/ccdc.csd.cc1zj31
Deposited on	29/03/2018
Crystallographer(s)	
Crystallographer	Andrew Bond
Affiliation	University of Cambridge

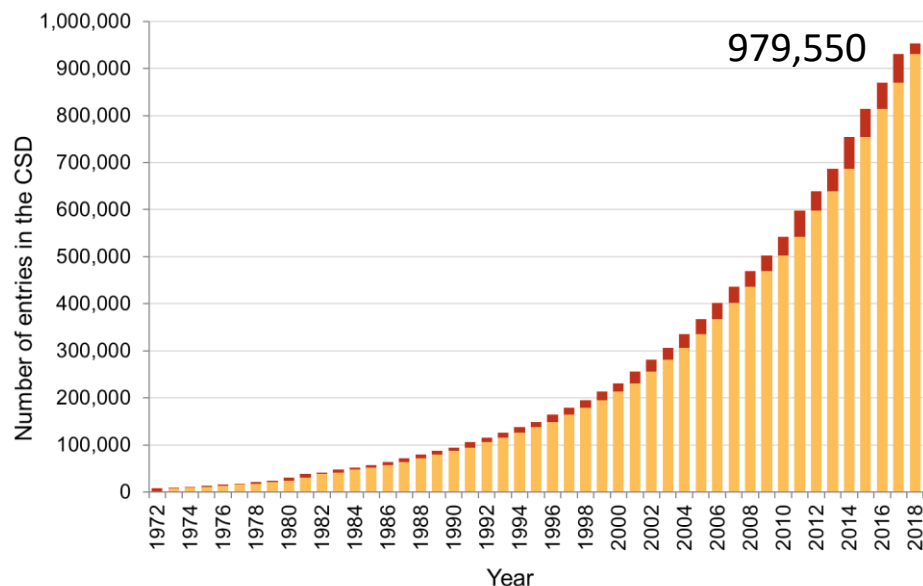
Associated publications



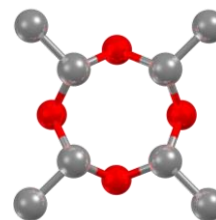
Alex J. Plajer, Annie L. Colebatch, Markus Enders, Álvaro García-Romero, Andrew D. Bond, Raúl García-Romero, Dominic S. Wright, *Dalton Transactions*, 2018, 47, 7036, DOI: 10.1039/C8DT01332B



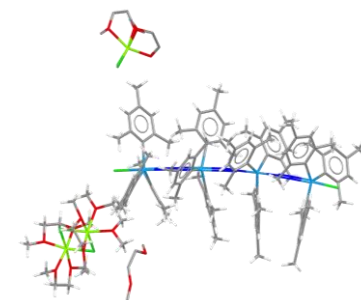
The Cambridge Structural Database



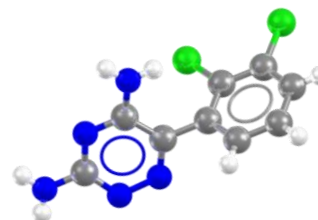
- ❑ 970,000+ small-molecule crystal structures
- ❑ Over 80,000 datasets deposited annually
- ❑ Structures available for anyone to download
- ❑ Links to over 1,000 journals
- ❑ Enriched and annotated by experts
- ❑ Access to data and knowledge



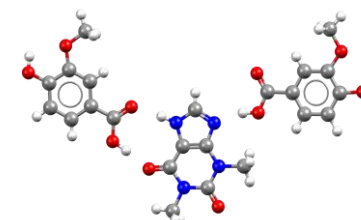
One of 1st – METALD
Metaldehyde – published 1936



250,000th – IBEZUK
Conducting metal-dinitrogen polymer



500,000th – EFEMUX01
Lamotrigine – an anti-convulsant drug



750,000th – ZOYBIA
Co-crystal of vanillic acid and theophylline

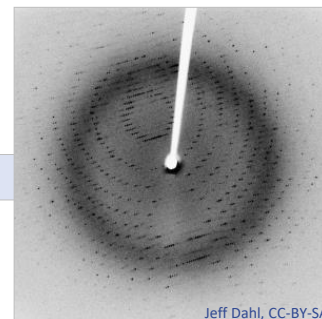
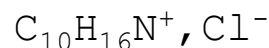


1,000,000th structure expected in
2019 #CSD1Million



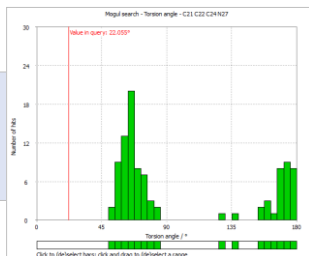
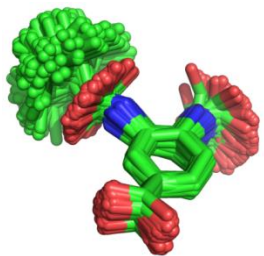
From Data to Knowledge

Experimental Data

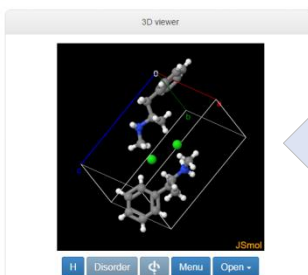


```
_diffrn_ambient_temperature 90 (2)
_diffrn_radiation_type MoK\alpha
_diffrn_radiation_wavelength 0.71073
_diffrn_radiation_monochromator graphite
_diffrn_measurement_device_type 'Bruker APEX CCD area-detector'
_diffrn_measurement_method '\psi and \omega'
_diffrn_detector_area_resol_mean 512
_diffrn_refine_number 5552

loop
  _atom_site_type_symbol
  _atom_site_label
  _atom_site_fract_x
  _atom_site_fract_y
  _atom_site_fract_z
  _atom_site_U_iso_or_equiv
  _atom_site_occupancy
  _atom_site_disorder_group
  _atom_site_refinement_flags
  _atom_site_symmetry_multiplicity
  _atom_site_disorder_assembly
  _atom_site_disorder_group
C1 C11 0.23185(8) 0.78305(9) 0.55674(6) 0.02219(16) Dani d U 1 1 . .
N H1 0.8031(3) 0.4811(3) 0.5363(2) 0.0172(4) Dani d U 1 1 . .
C C1 0.4936(4) 0.7587(6) 0.4357(2) 0.0245(8) Dani d U 1 1 . .
C C2 0.7510(5) 0.8922(5) 0.7083(3) 0.0256(6) Dani d U 1 1 . .
C C3 0.7409(4) 0.6944(4) 0.6644(3) 0.0187(6) Dani d U 1 1 . .
C C4 0.8700(4) 0.6537(4) 0.7485(3) 0.0234(6) Dani d U 1 1 . .
```



MIYCLUT (+)-Methamphetamine hydrochloride
Spacegroup P21, Cell a 7.1022(11)Å b 7.2949(11)Å c 10.8121(17)Å α 90°
β 97.295(4)° γ 90°



organic compounds

Acta Crystallographica Section E
Structure Reports
Online
ISSN 1600-5368

Redetermination of (+)-methamphetamine hydrochloride at 90 K

Trick Hakey, Wayne Ouellette, Jon Zubieta and Timothy ...

Department of Chemistry, Syracuse University, Syracuse, New York 13244, USA
E-mail: oellette@atlas.syr.edu

Received 20 March 2008; accepted 22 April 2008

Key indicators: single-crystal X-ray study; T = 90 K; mean σ(C–C) = 0.004 Å; R factor = 0.052; wR factor = 0.118; data-to-parameter ratio = 15.6.

The title crystal structure of the blocky name: N-methyl-1-phenylpropan-2-aminium chloride), C₁₀H₁₆N⁺Cl⁻, was originally determined by Simon, Ioske & Torok [Acta Pharm. Mater. (1987), 63, 334–330] and Van, Kim, B. Min, H. Hwang ...

Data collection: SMART (Bruker, 2002); data reduction: SAINT (Bruker, 2002); structure solution: SHELXS (Sheldrick, 2002); structure refinement: SHELXL (Sheldrick, 2002); molecular graphics: ORTEP-3 (Johnson, 1999); data collection: SMART (Bruker, 2002); data reduction: SAINT (Bruker, 2002); structure solution: SHELXS (Sheldrick, 2002); structure refinement: SHELXL (Sheldrick, 2002); molecular graphics: ORTEP-3 (Johnson, 1999); publication material: CIF2000 (IUCr, 2000); software: Bruker AXS (2002).

Refinement	
R ² = 0.052	wR ² = 0.117
S = 1.05	
2720 reflections	
174 parameters	
1 restraint	

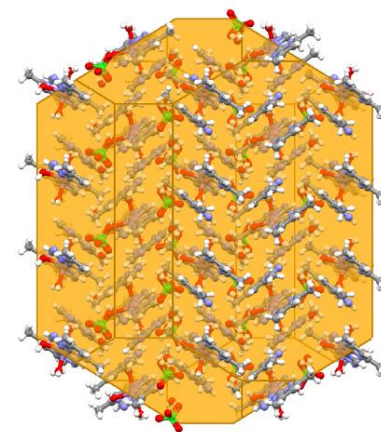
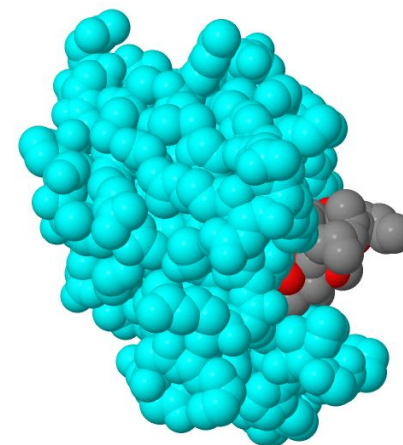
Table 1 Hydrogen-bond geometry (Å, °)		
D—H...A	D—H	H...A
N1—H11D...C11 ⁱ	0.93 (4)	2.34 (4)
N1—H11E...C11 ⁱ	0.90 (3)	2.25 (4)
Symmetry codes: (i) −x + 1, y − 1/2, −z + 1/2 (iii)		

Structural Knowledge



Structural Chemistry Insights

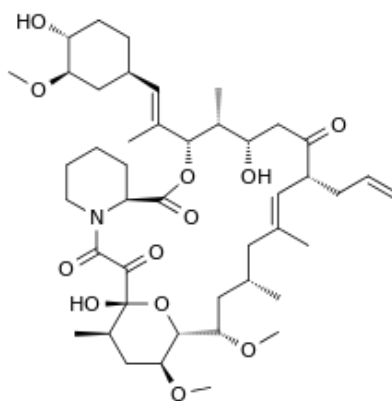
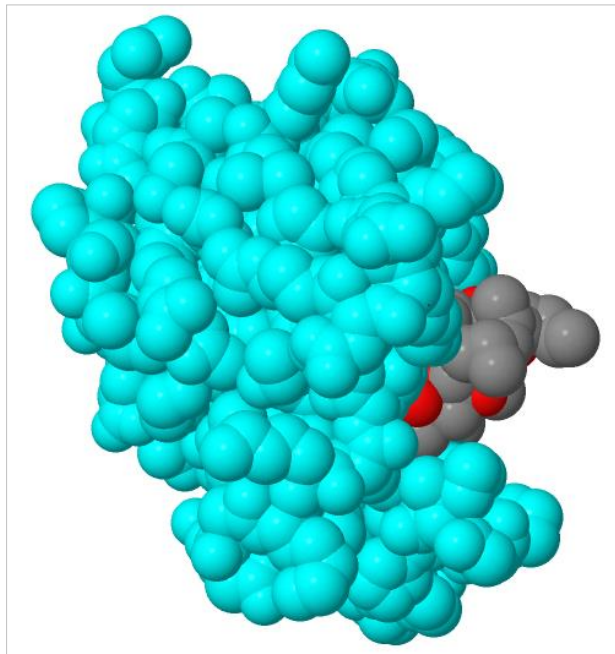
- **Crystal structure data gives insights into:**
 - Molecular dimensions and shape
 - Molecular interactions
 - Solid form properties
- **Applicable to various domains including:**
 - Drug design and development
 - Design of new materials
 - Crystal engineering
 - Structure validation





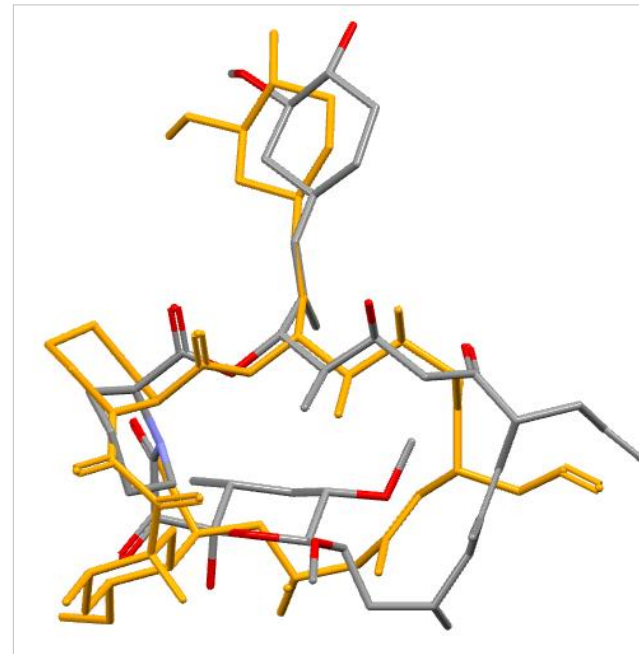
Crystal Structure Knowledge Helps Drug Design

PDB 1BKF complexed with PDB Ligand FK5



Tacrolimus: An immunosuppressive drug. Also used in the treatment of skin conditions.

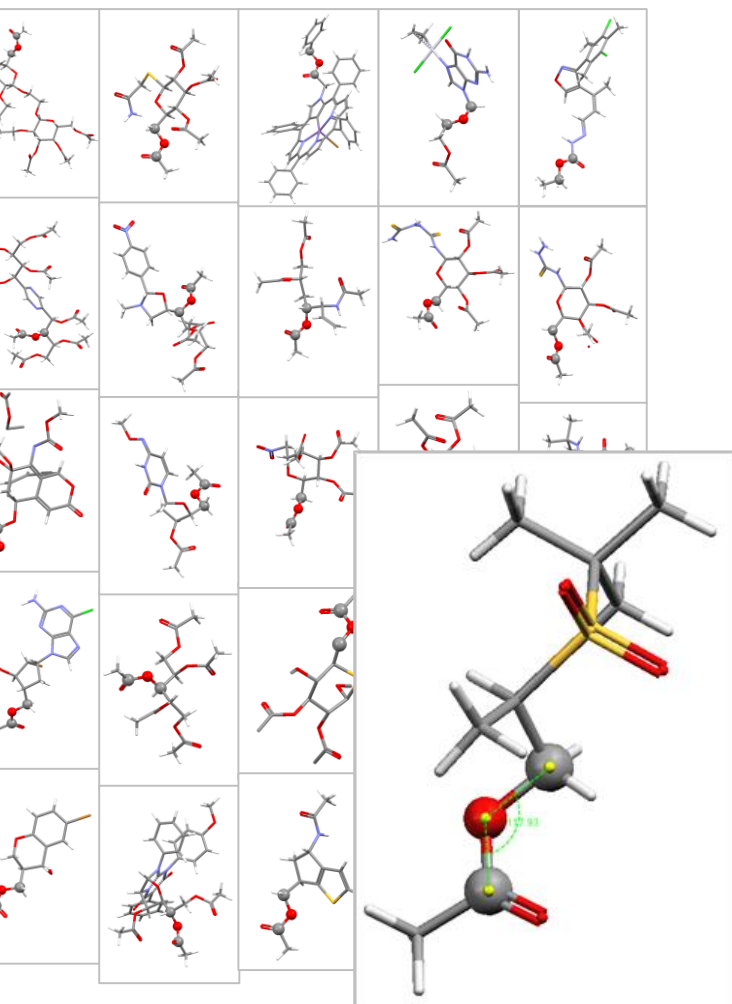
CSD FINWEE10 overlayed on PDB Ligand FK5



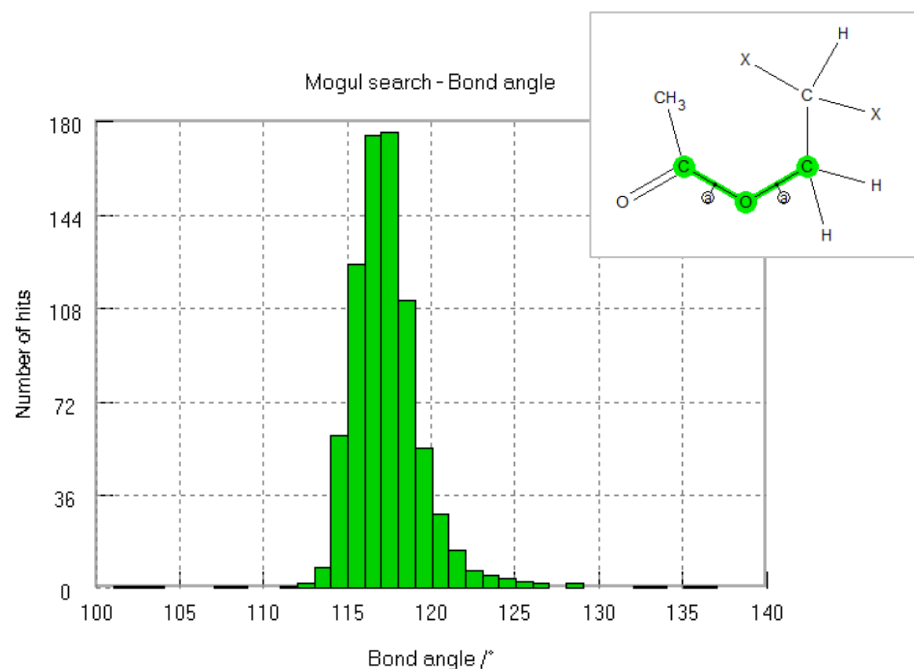
Understanding factors that influence the shape of molecules helps identify better drug candidates



Mogul A Knowledge Base of Molecular Geometries

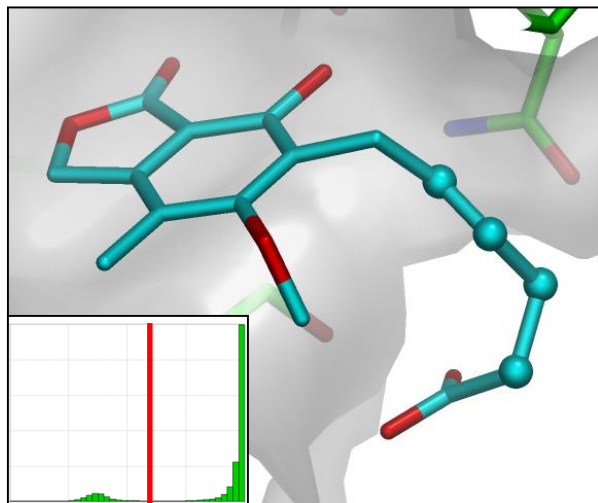


Millions of experimentally observed bond lengths, valence angles and torsions are combined into distributions that indicate preferred geometries of structural features



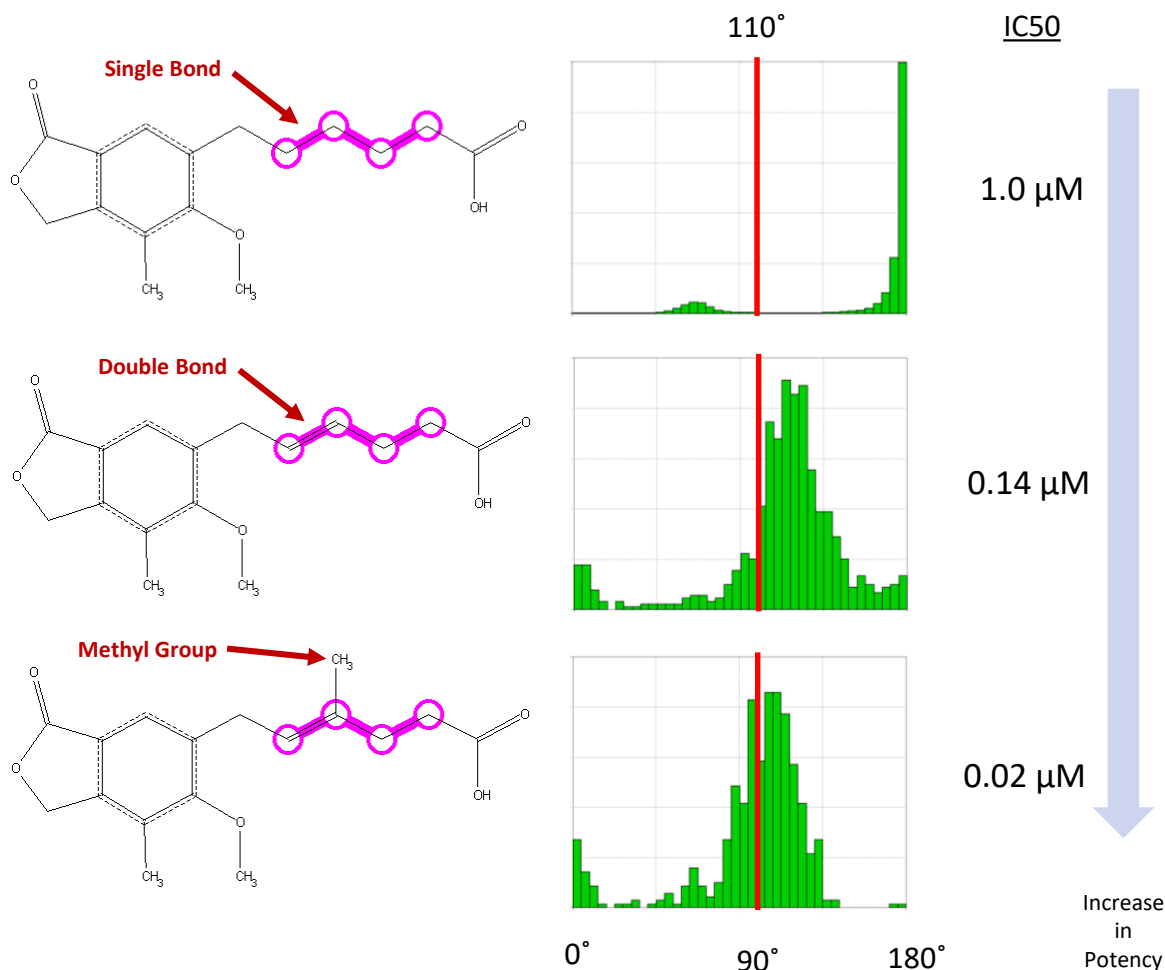


Knowledge-driven Conformer Design



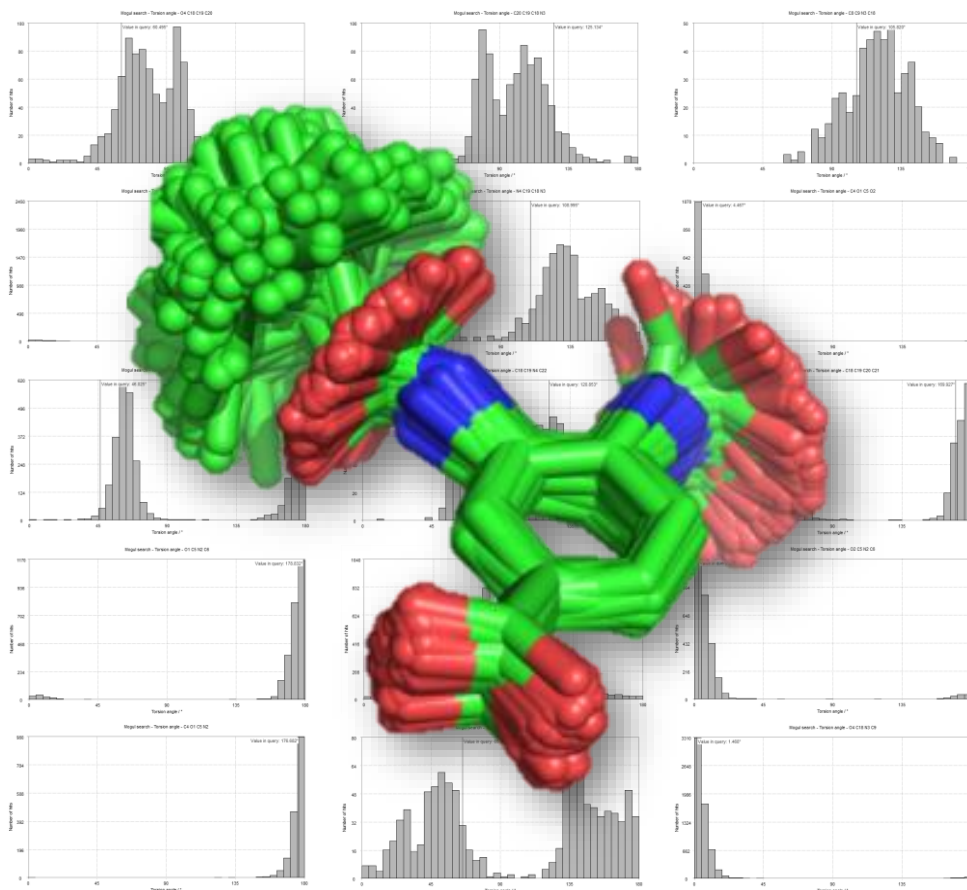
The immunosuppressant mycophenolic acid binds to inosine monophosphate dehydrogenase with an affinity of 1.0 μM .

The CSD reveals that the highlighted torsion at 110° is unfavourable.





Knowledge-based Conformer Generation



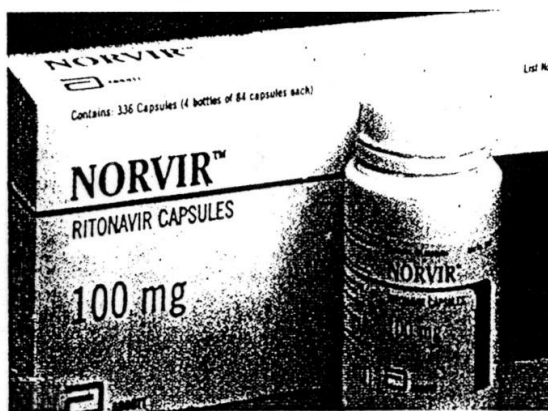
- Bond, angle and torsion distributions derived from experimental structures are used to produce realistic ensembles of low energy structures.
- The CSD Conformer Generator can be used to minimise molecular conformations and generate diverse conformer subsets based on experimental data
- Applicable to drug discovery and drug development



Mitigating Risk in Drug Development

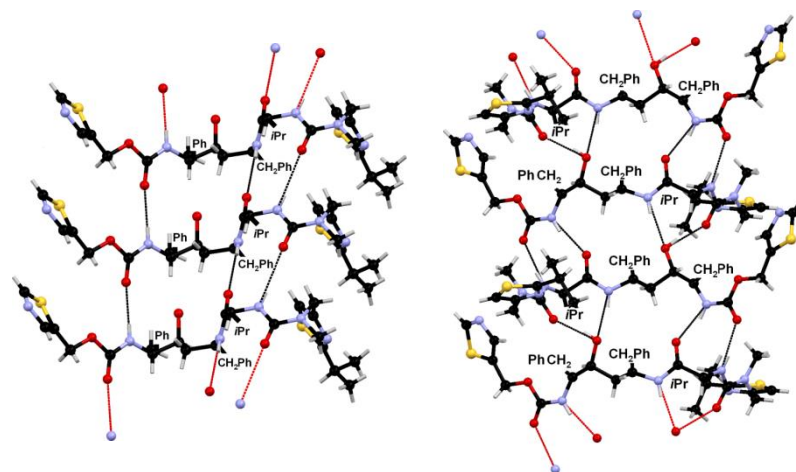
Manufacturing problems hit Abbott's HIV drug ritonavir

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.



Capsules unlikely to be available from mid-August

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

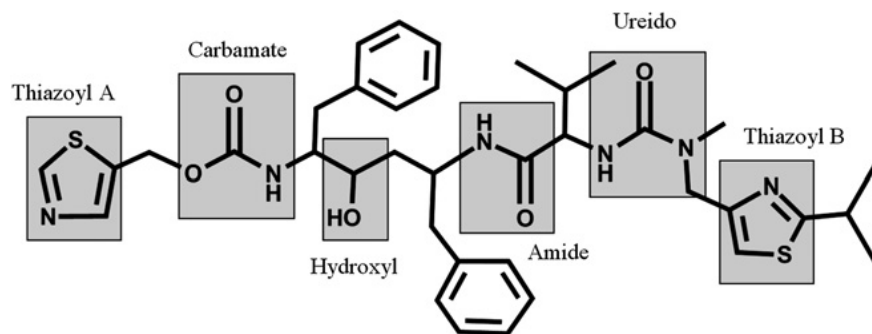


Different crystal forms, different interactions, different solubility, different stability.

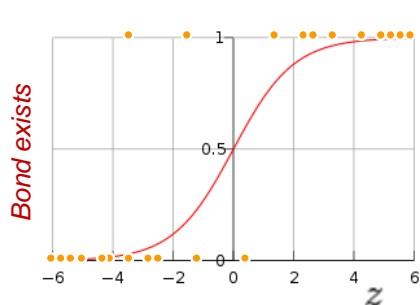
Knowing the likelihood of specific molecular interactions occurring helps assess the risk of undesirable crystal formation



Hydrogen Bond Propensity: Ritonavir



Predictive analytics is used to identify feasible and unusual crystal packings based on information from known crystal structures of molecules similar to the target.



$$\text{Propensity } \Pi \text{ (H-bond exists)} = \frac{1}{1 + e^{-z}}$$

$$Z = \beta_0 + \beta_1 \cdot \text{DG1} + \beta_2 \cdot \text{AG1} + \beta_3 \cdot \text{DSD} + \dots$$

DG1, AG1, DSD: Explanatory variables

Logistic Regression

Table 1 Propensity predictions for potential donor–acceptor combinations in ritonavir (as labelled in Fig. 1), and observed hydrogen bonds in either polymorphic form

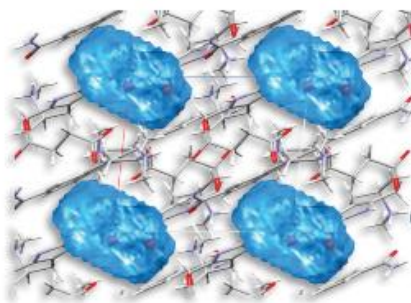
Donor	Acceptor	π	\pm^a	Form I	Form II
amide	carbamate	0.618	0.094	X	X
amide	hydroxyl	0.551	0.052	X	✓
carbamate	carbamate	0.538	0.090	✓	X
hydroxyl	carbamate	0.537	0.090	X	X
amide	amide	0.501	0.055	✓	X
amide	ureido	0.499	0.072	X	X
carbamate	hydroxyl	0.470	0.078	X	X
hydroxyl	hydroxyl	0.469	0.037	X	X
carbamate	amide	0.420	0.083	X	✓
hydroxyl	amide	0.419	0.045	X	X
carbamate	ureido	0.418	0.088	X	X
hydroxyl	ureido	0.417	0.058	X	✓
ureido	carbamate	0.319	0.086	X	✓
ureido	hydroxyl	0.263	0.041	X	X
ureido	amide	0.225	0.040	X	X
ureido	ureido	0.224	0.044	✓	X
amide	thiazoyl a	0.152	0.054	X	X
amide	thiazoyl b	0.142	0.050	X	X
carbamate	thiazoyl a	0.115	0.044	X	X
hydroxyl	thiazoyl a	0.114	0.039	✓	X
carbamate	thiazoyl b	0.107	0.041	X	X
hydroxyl	thiazoyl b	0.106	0.036	X	X

^a The error bars of the coefficient value: the value falls within this range at the 95% confidence level, based on a χ^2 distribution.



CSD-Materials

Informatics-based solutions that aid in the understanding and prediction of solid form stability and properties

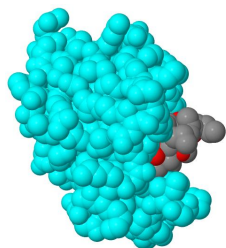


- Discover preferred interactions and engineer changes to satisfy these requirements using *Full Interaction Maps*
- Interpret crystal packing and compare with CSD data using powerful *packing feature, similarity and motif searches*, and *hydrogen bond propensity analysis*
- Understand the effects of hydration on your lattice with the *hydrate analyser*
- Explore the structures of potential co-crystals using the *molecular complementarity tool*

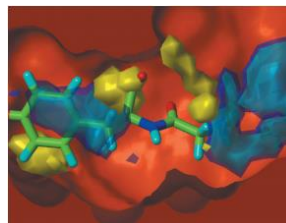
Developed in partnership with Industry through CCDC's Crystal Form Consortium



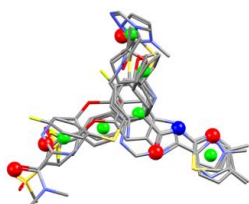
CSD-Discovery



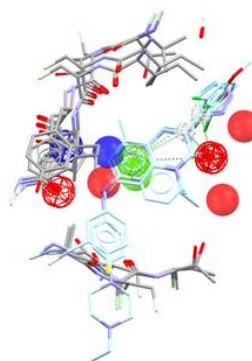
GOLD: Protein-ligand docking - virtual screening, lead optimisation and binding mode prediction



SuperStar: Knowledge-based prediction of intermolecular interactions based on data from small molecule crystal structures



Ligand Overlay: applying structural knowledge to identify common binding modes, interactions and geometries of structurally diverse ligands



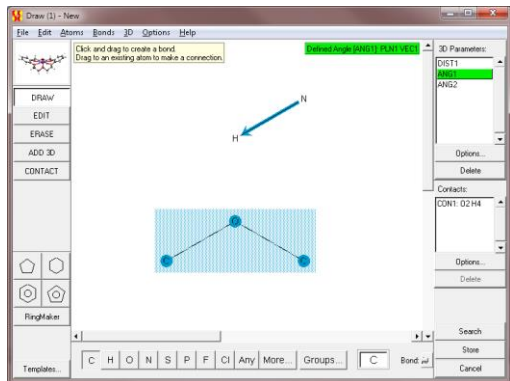
CSD-CrossMiner: Fast and flexible pharmacophore searching across the CSD and the PDB for lead optimisation and scaffold hopping



CSD-System



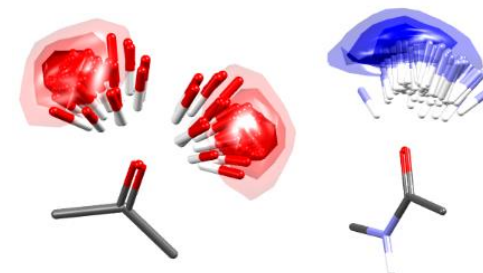
ConQuest: Advanced 3D searching



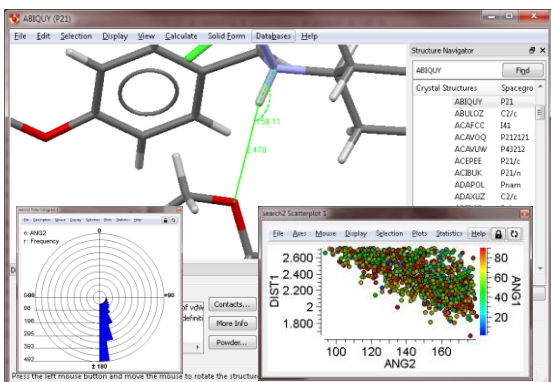
WebCSD: On-line portal to the CSD



IsoStar: Molecular interaction analysis



Mercury: Visualisation & data analysis



CSD Python API: Custom search and analysis

```
import sys
from ccdc import conformer
from ccdc import io

args = parser.parse_args()

mol_reader = io.MoleculeReader(args.inmolfn)
engine = conformer.GeometryAnalyser()

molecules = []
min_unusual_torsions = sys.maxint
for (idx, molecule) in enumerate(mol_reader):
    molecule.standardise_aromatic_bonds()
    molecule.standardise_delocalised_bonds()

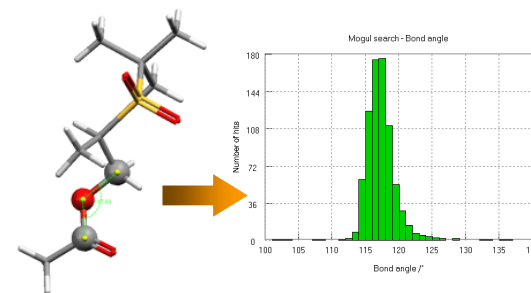
    # Do the analysis
    geometry_analysed_molecule = engine.analyse_molecule(molecule)

    # Count number of unusual torsions
    molecule.unusual_torsions = []
    for t in geometry_analysed_molecule.analysed_torsions:
        if t.unusual and t.enough_hits:
            num_unusual_torsions = len(molecule.unusual_torsions)
            molecule.num_unusual_torsions = num_unusual_torsions
            molecules.append(molecule)

    if num_unusual_torsions < min_unusual_torsions:
        min_unusual_torsions = num_unusual_torsions
```



Mogul: Molecular geometry analysis





CCDC Knowledge-based Software Solutions

CSD-System: Crystallographers, structural chemists, educators

Find, analyse and communicate crystal structures

CSD-Discovery

Medicinal chemists, computational chemists, structural biologists

Protein and ligand-based design of new molecules

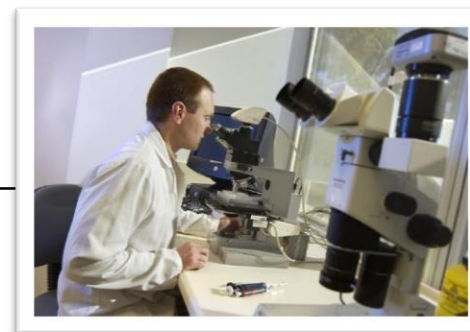
CSD-Materials

Solid form experimentalists, crystallization scientists

Behaviour and properties of new materials



CSD-Enterprise
All CCDC application software
(available to all Academics)



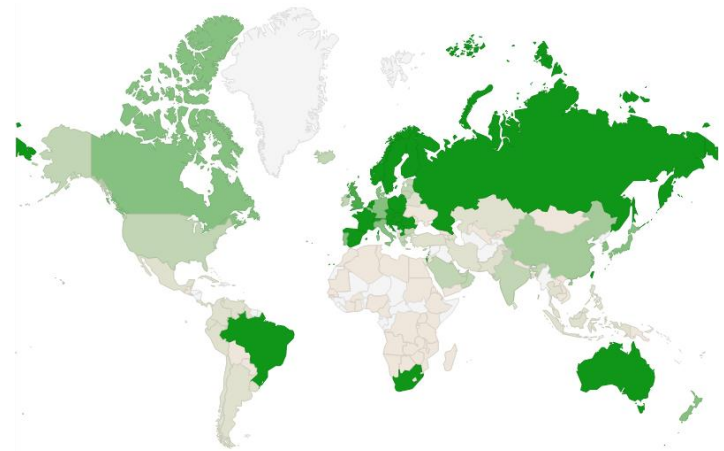


Academic Access to Value-added Services

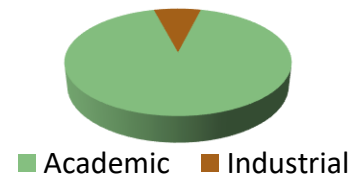
- **Country-wide licences**
 - via National Affiliated Centres
- **Campus-wide licences**
 - often through University Libraries
- **Individual Researchers and Groups**
 - cost can be included in grant applications
- **CCDC Subsidy**
 - Supporting research in developing countries via the [FAIRE](#) programme

24 countries with academic country-wide licences

Colour-coding reflects nature of academic arrangements

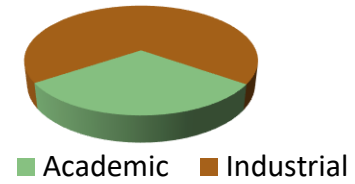


User Community



Number of customers, 2018

Revenue



CSD and Software, 2018



Chemistry Department Access to CSD-Enterprise

- **Download CSD onto your machine** (Linux, Windows, OS X) following these instructions: <https://www.ch.cam.ac.uk/computing/software/cambridge-structural-database-system>
- CSD-System Components **installed on MCS machines** in the Chemistry Library and around the University:
 - CSD Conquest
 - CSD DASH
 - CSD 2017 encipher
 - CSD 2017 GOLD
 - CSD 2017 Hermes
 - CSD 2017 IsoStar
 - CSD 2017 IsoStar Server
 - CSD 2017 Mercury
 - CSD 2017 Mogul
 - CSD 2017 PreQuest
 - CSD 2017 Python API

The screenshot shows the WebCSD interface. At the top, there is the Cambridge Crystallographic Data Centre logo and the text 'WebCSD'. Below this are search tabs: 'Simple Search', 'Structure Search', 'Unit Cell Search', and 'Formula Search'. A search bar contains the text 'Chemical structure searching'. Below the search bar, there is a prompt: 'Please draw your diagram or add a SMARTS string in the 'advanced' section below.' A drawing toolbar is visible, and the main area shows a chemical structure of a silicon atom bonded to two methyl groups and two 2-methylpyridine rings. The silicon atom is labeled 'Si' and has a bond to a group labeled '?R2'. At the bottom, there are match condition options: 'Exact', 'Substructure', and 'Similarity'. A 'Search' button and a 'Clear' button are also present. The dotmatics logo is visible in the bottom right corner of the drawing area.

<https://webcsd.ccdc.cam.ac.uk/>



Sharing Crystallographic Data

Principles and Practices



CSD-Community Services

- Free community services

- Data deposition
- Data validation
- Data archiving
- Data access

CCDC identifier(s)

Compound name

DOI

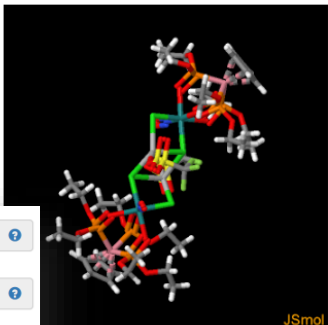
Authors

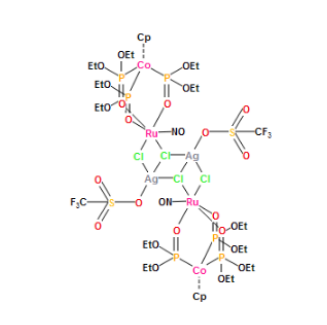
Journal

Publication details

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

LIMSUV : bis((μ-Cloro)-(μ-chloro)-hexakis(μ-diethyl phosphonato)-bis((η⁵-cyclopentadienyl)-(trifluoromethanesulfonate)-(nitrosyl))-di-cobalt-di-ruthenium(II)-di-silver
Space Group: P21/n, Cell: a 9.5621(6)Å, b 36.808(2)Å, c 19.0222(11)Å, α 90° β 91.1430(10)° γ 90°

3D viewer  JSmol

Chemical diagram 

Style Labels Packing Measure
d Sticks No Labels None None

[View group symbols key](#)

Additional CCDC details

Citation - Xiao-Yi Yi, T.C.H.Lam, Yiu-Keung Sau, Qian-Feng Zhang, I.D.Williams, Wa-Hung Leung, CCDC
6: Experimental Crystal Structure Determination, 2007, DOI: [10.5517/ccq8y74](https://doi.org/10.5517/ccq8y74)
dated on: 15/10/2007

Associated publications

Xiao-Yi Yi, T.C.H.Lam, Yiu-Keung Sau, Qian-Feng Zhang, I.D.Williams, Wa-Hung Leung, *Inorganic Chemistry*, 2007, 46, 7193, DOI: [10.1021/ic7007683](https://doi.org/10.1021/ic7007683)

0.044(3) Uani 0.50 1 d PDU A 1
.0327(11) Uani 0.50 1 d PDU A 1
.029(4) Uani 0.50 1 d PDU A 1
0.031(4) Uani 0.50 1 d PDU A 1
o 0.50 1 calc PR A 1
0.027(4) Uani 0.50 1 d PDU A 1
0.029(4) Uani 0.50 1 d PDU A 1
0.0586(15) Uani 0.50 1 d PDU A 1
0.0400(10) Uani 0.50 1 d PDU A 1

Funder Research Data Sharing Policies

“Publicly funded research data are a public good, (...), which should be made openly available with as few restrictions as possible in a timely and responsible manner (...).”

UK Research
and Innovation

“Research data that supports publications must be stored for 10 years.”



“EPSRC has the strictest policy on research data sharing...”



- All publications ... should have a statement describing how to access underlying data
- Data should be stored for at least 10 years or for 10 years from the last request for access to the data
- All data supporting research publications should be cited with the use of persistent links, for example DOIs (Digital Object Identifiers)...
- Research data should be accompanied by metadata ... to allow the discovery of data.

<https://www.data.cam.ac.uk/funders/epsrc-funded-researchers>



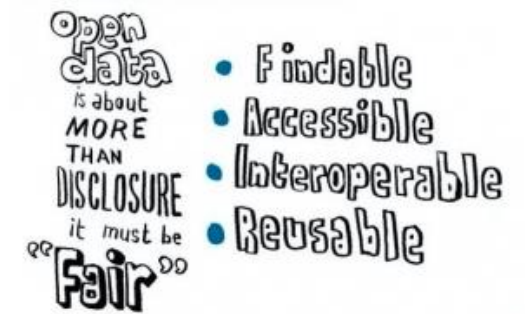
FAIR Data Principles

Comment | [OPEN](#)

The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier [...] Barend Mons

Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* 3:160018 doi: 10.1038/sdata.2016.18 (2016).



DATA SHOULD BE

Findable

Interoperable

Accessible

Reusable

BY HUMANS AND MACHINES



<http://www.datafairport.org/>



The Future of Research Communications and e-Scholarship

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



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**



CCDC Deposition Services

- 1 Login
- 2 Upload
- 3 Check Syntax
- 4 Validation
- 5 Add Publication
- 6 Enhance Data
- 7 Review
- 8 Submit

CIF deposition and validation service

First name(s)

Last name(s)

Your email address

Your ORCID ID Create or Connect your ORCID ID

Additional email addresses

Institution (e.g. University/Company)

Deposition number(s) for revision

CIF/HKL/RES/FCF/Word/ZIP files

- structure01.cif 22.82 KB
- structure02.cif 146.20 KB

Details Remember my details

Options I wish to run the IUCr *checkCIF/PLATON* service on my data



CCDC Deposition Services

- 1 Login
- 2 Upload
- 3 Check Syntax
- 4 Validation
- 5 Add Publication
- 6 Enhance Data
- 7 Review
- 8 Submit

CIF deposition Check Syntax

The files highlighted in red in the left-hand column contain errors that need fixing before proceeding.

Please click on any red file names in the left-hand column, make the appropriate edits and then click the 'Save & Recheck File' button before proceeding to the next step.

For more information on how to fix errors please see our [correcting CIFs](#) page.

Pick file to edit

structure01.cif

structure02.cif

File contents structure01.cif

```
30 data_I
31 _audit_creation_method SHELXL-97
32 _chemical_name_systematic
33 ;
34 {5-[(7-chloroquinolinium-4-yl)amino]-2-hydroxybenzyl}diethylammonium dichloride
35 dihydrate
36 ;
37 _chemical_name_common Amodiaquine dihydrochloride dihydrate'
38 _chemical_formula_moiety 'C20 H24 Cl N3 O 2+, 2(Cl -), 2(H2 O)'
39 _chemical_formula_sum C20 H28 Cl3 N3 O3'
40 _chemical_formula_iupac 'C20 H24 Cl N3 O 2+, 2Cl -, 2H2 O'
41 _chemical_formula_weight 464.80
42 _chemical_melting_point ?
43 _symmetry_cell_setting monoclinic
44 _symmetry_space_group_name_H-M 'P 21/c'
45 _symmetry_space_group_name_Hall '-P 2ybc'
46 loop_
47 _symmetry_equiv_pos_as_xyz
48 'x, y, z'
49 '-x, y+1/2, -z+1/2'
50 '-x, -y, -z'
51 'x, -y+1/2, z+1/2'
52 _cell_length_a 7.76220(10)
53 _cell_length_b 26.8709(4)
54 _cell_length_c 10.7805(2)
55 _cell_angle_alpha 90.00
56 _cell_angle_beta 92.7040(10)
57 _cell_angle_gamma 90.00
58 _cell_volume 2230.91(6)
59 _cell_formula_units_Z 4
```

← Go Back

Save & Recheck File

Proceed to Next Step →

Error 44 No terminating () quote

No Structure Factor data

Structure Factor data are a
you are unable to embed S

If in exceptional circumstan
Next Step', however you m
embedded into your deposi

Reason why your deposition c



Crystallographic Information File: CIF

Acta Cryst. (1991). A47, 655–685

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

(Received 8 April 1991; accepted 28 June 1991)

- Data items semantically defined by CIF dictionaries (vocabularies)
 - crystallisation details
 - instrument details
 - software packages and parameters
 - quality metrics
 - publication details

- 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 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 Integrity: checkCIF



checkCIF

A service of the
International Union of Crystallography

checkCIF reports on the consistency and integrity of crystal structure determinations reported in CIF format.

Please upload your CIF using the form below.

File name:

No file chosen

Select form of checkCIF report

HTML PDF

Select validation type

Full validation of CIF and structure factors
 Validation of CIF only (no structure factors)

Output Validation Response Form

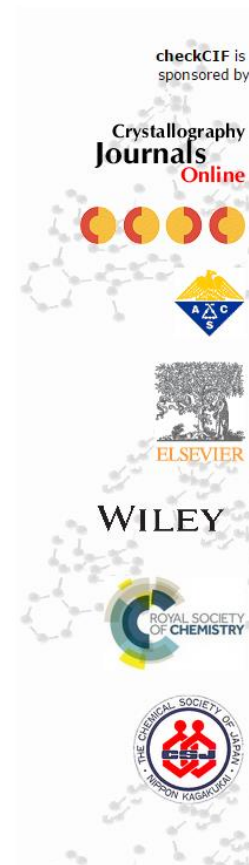
Level A alerts only
 Level A and B alerts
 Level A, B and C alerts
 None



- Checks consistency and integrity of the data
- Generates alerts that should either be corrected or explained

Level A	Most likely a serious problem - resolve or explain
Level B	A potentially serious problem, consider carefully
Level C	Check. Ensure it is not caused by an omission or oversight
Level G	General information/check it is not something unexpected

Much of checkCIF based on components of PLATON developed by Ton Spek, Utrecht University





CCDC Deposition Services

- 1 Login
- 2 Upload
- 3 Check Syntax
- 4 Validation
- 5 Add Publication
- 6 Enhance Data
- 7 Review
- 8 Submit

CIF dep Check Syntax

Validation

View reports on the consistency and integrity of your structures

Datablock: sa2906a

Bond precision: C-C = 0.0118 Å Wavelength=0.71073
Cell: a=6.991(5) b=10.778(5) c=15.575(5)
alpha=90 beta=90 gamma=90
Temperature: 293 K

	Calculated	Reported
Volume	1173.6(11)	1174(2)
Space group	P 21 21 21	P212121
Hall group	F 2ac 2ab	?
Moiety formula	C8 H10 Cu N4 O4, H2 O	C8 H10 Cu N4 O4,
Sum formula	C8 H12 Cu N4 O5	C8 H12 Cu N4 O5
Mr	307.77	307.77
Dx, g cm-3	1.742	1.742
Z	4	4
Mu (mm-1)	1.882	1.882
F000	628.0	680.0
FO00'	629.54	
h, k, lmax	9, 14, 21	9, 13, 20
Nref	3036[1759]	2901
Tmin, Tmax		
Tmin'		
Correction method=	Not given	
Data completeness=	1.65/0.96	Theta(max)= 28.660
R(reflections)=	0.0690(1863)	wR2(reflections)= 0.1805(2901)
S =	0.913	Npar= 169

The following ALERTS were generated. Each ALERT has the format
test-name_ALERT alert-type alert-level.
Click on the hyperlinks for more details of the test.

Alert level A
[EXPT005_ALERT_1_A](#) _exptl_crystal_description is missing
Crystal habit description.
The following tests will not be performed.
CRYSR_01
[DIFF003_ALERT_1_A](#) _diffn_measurement_device_type is missing
Diffractometer: make and type. Replaces _diffn_measurement_device_type.
[PLAT183_ALERT_1_A](#) Missing _cell_measurement_reflns_used Value Please Do!
[PLAT184_ALERT_1_A](#) Missing _cell_measurement_theta_min Value Please Do!
[PLAT185_ALERT_1_A](#) Missing _cell_measurement_theta_max Value Please Do!

Alert level B
[PLAT420_ALERT_2_B](#) D-H Without Acceptor O4 --H2B ..

IUCr checkCIF ?

Unit cell check ?

IUCr checkCIF Response
Please enter your response here for structure02.cif / data_sa2906a.

Level A

- [EXPT005_exptl_crystal_description](#) is missing
- [DIFF003_diffn_measurement_device_type](#) is missing
- [PLAT183](#) Missing _cell_measurement_reflns_used Value Please Do!
- [PLAT184](#) Missing _cell_measurement_theta_min Value Please Do!
- [PLAT185](#) Missing _cell_measurement_theta_max Value Please Do!

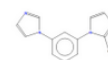
Level B

- [PLAT420](#) D-H Without Acceptor O4 --H2B .. Please Check

Level C

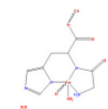
Save Close

FAMCUU



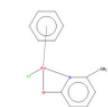
Deposition Number(s): 1501474
Space Group: P 2₁ 2₁ 2₁ (19)
Cell: a 6.9209(2)Å b 10.5769(2)Å c 15.4237(4)Å, α 90° β 90° γ 90°

FINXAD



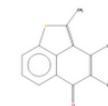
Deposition Number(s): 912870
Space Group: P 2₁ 2₁ 2₁ (19)
Cell: a 6.991(5)Å b 10.778(5)Å c 15.575(5)Å, α 90.00° β 90.00° γ 90.00°

GIDXUM



Deposition Number(s): 1167512
Space Group: P 2₁ 2₁ 2₁ (19)
Cell: a 6.844(2)Å b 10.975(5)Å c 15.363(5)Å, α 90° β 90° γ 90°

KELCAH



Deposition Number(s): 886135
Space Group: P 2₁ 2₁ 2₁ (19)
Cell: a 6.953(1)Å b 10.665(2)Å c 15.414(3)Å, α 90° β 90° γ 90°



CCDC Deposition Services

- 1 Login
- 2 Upload
- 3 Check Syntax
- 4 Validation
- 5 Add Publication
- 6 Enhance Data
- 7 Review
- 8 Submit

CIF deposition

Check Syntax

Validation

Add Publication

Please check and add/update the publication details shown below.

If you don't know the full publication details then please provide the current list of authors for

Authors ? *

Nicola Lavery, Arthur Smith

Journal name ?

Volume ?

Volume

Year ?

Year

Page ?

Page

Publication DOI ?

E.g. 10.14469/hpc/2300

Additional information ?

If you do not intend to publish your data in the scientific literature immediately through the Cambridge Structural Database (CSD) (ICSD) then please click the 'Publish in a Database' button below. Inorganic data should be published in the CSD as a [CSD Communication](#). Inorganic data should be published in the CSD as a [CSD Communication](#).

Publish in a Database

Error 44 No terminating (') quote

Add Crystallographer Details

Please add the details of the main crystallographer associated with the data below. The email address will be used to notify the crystallographer about this deposition. The name, affiliation, country and if appropriate ORCID iD of the crystallographer may be displayed to users alongside the data.

CSD Communications

The CCDC allows you to publish data directly through the CSD as a *CSD Communication*. Over 5,000 structures were published in this way in 2016 making *CSD Communications* the No. 1 place to publish crystal structures.



Recognition

A citable DOI allows you to receive credit for your structures and add the data to your ORCID record.

If you would like to reference a CSD Communication then we would recommend using the following style of citation

"D.S.Cati, H.Stoeckli-Evans, CCDC 227635: *CSD Communication*, 2004, DOI: 10.5517/cc7mw2s"



Discoverability

Automatic linking via CCDC DOI from third-party repositories, such as DataCite, the Web of Science Data Citation Index and may also be linked from ChemSpider, PubChem and PDB Chemical Component Dictionary



Identifiers and Data Citation



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](https://doi.org/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>



Linking Researchers, Datasets and Articles

The Cambridge Crystallographic Data Centre

CSD Entry: GIGMIV

GIGMIV : bis(1,4,7,10-tetraoxacyclododecane)-lithium 3-methyl-1,2-benzodiphosphol-1-ide
Space Group: $P \bar{1} (2)$, Cell: a 8.4645(4)Å b 12.2344(6)Å c 14.0658(7)Å, α 95.494(3)° β 95.567(4)° γ 110.021(3)°


Database Identifier	Deposition Number
<input checked="" type="checkbox"/> GIGMIV	1852044
<input checked="" type="checkbox"/> GIGMOB	1852045

[Download](#)


Additional details

Deposition Number	1852044
Data Citation	Lily S. H. Dixon, Schirin Hanf, Jessica E. Waters, Andrew D. Bond, Dominic S. Wright CCDC 1852044: Experimental Crystal Structure Determination, 2018, DOI: 10.5517/ccdc.csd.cc2056cs
Deposited on	27/06/2018

Crystallographer(s)

Crystallographer	Andrew Bond 
Affiliation	University of Cambridge

Associated publications

 Lily S. H. Dixon, Schirin Hanf, Jessica E. Waters, Andrew D. Bond, Dominic S. Wright, *Organometallics*, 2018, 37, 4465, DOI: [10.1021/acs.organomet.8b00480](https://doi.org/10.1021/acs.organomet.8b00480)

Data Citation with Dataset DOI

Link to Crystallographer's ORCID Record

Link to article based on article DOI



CCDC Deposition Services

- 1 Login
- 2 Upload
- 3 Check Syntax
- 4 Validation
- 5 Add Publication
- 6 Enhance Data
- 7 Review
- 8 Submit

CIF deposition

Check Syntax

Validation

Enhance Data

Add Publication

Please check and add/update the publication details shown below.

If you don't know the full publication details then please provide the following information:

Authors [?] * Nicola Lavery, Arthur Smith

Journal name [?]

Volume [?] Volume

Year [?] Year

Page [?] Page

Publication DOI [?] E.g. 10.14469/hpc/2300

Additional information [?]

If you do not intend to publish your CIF file immediately through the Cambridge Crystallographic Data Centre (CCDC) then please click the 'Publish in the CSD as a [CSD Communication](#)' button below.

Publish in a Database

Error 44

Pick a structure to edit

structure01.cif

data_1

structure02.cif

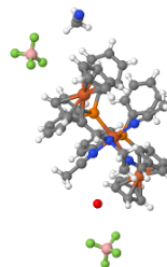
data_sa2906c

data_sa2906a

data_sa2906b

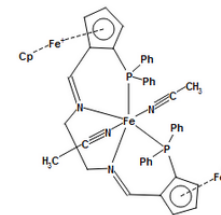
data_sa2906g

3D viewer



JSmol

Chemical diagram [?]



NC-CH₃ BF₄⁻ H₂O

data_sa2906c

```
21 # For further information on the CCDC, data deposition and
22 # data retrieval see:
23 # www.ccdc.cam.ac.uk
24 #
25 # Bona fide researchers may freely download Mercury and enCIFer
26 # from this site to visualise CIF-encoded structures and
27 # to carry out CIF format checking respectively.
28 #
29 data_sa2906c
30
31
32 _audit_creation_method SHELXL-97
33 _chemical_name_systematic
34 ;
35 Bis(acetonitrile- $\kappa$ N)-((2Rp,2''Rp)-1,1''-
36 [1,2-ethanediybis(nitrilomethylidene)]bis[2-(diphenylphosphino)-ferrocene]-
37  $\kappa$ 4-N,N',P,P')-iron(ii) tetrafluoroborate acetonitrile water solvate 1/0.5
38 ;
39 _chemical_name_common '1704 MB209 at 200K'
40 _chemical_melting_point ?
41 _chemical_formula_moiety
42 'C52 H48 Fe3 N4 P2, 2(B F4), C2 H3 N, 0.5(H2 O)'
43 _chemical_formula_sum 'C54 H52 B2 F8 Fe3 N5 O0.50 P2'
44 _chemical_formula_weight 1182.12
45
46 loop_
47 _atom_type_symbol
```

Associated DOIs

Raw data DOI [?]

Data fields

Compound name [?]

Bis(acetonitrile- κ N)-((2R_p,2''R_p)-1,1''-[1,2-ethanediybis(nit

Synonyms/other names [?]

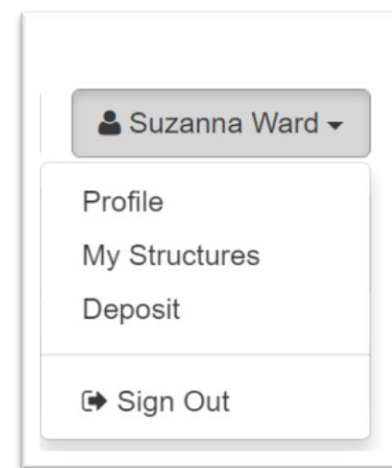
Crystal colour [?]

Crystal habit [?]



My Structures

- **Personalised deposition portal launched December 2016**
 - Ability to log on and view and retrieve depositions
 - Deposit and revise data
 - Edit and update basic information
 - Publish data directly as a *CSD Communication*
 - Share data with co-workers



My Structures

Search by CCDC Number

After depositing structures to the CCDC they may take a few minutes to appear in the table. If you are still unable to see your structures after you have received your CCDC numbers please contact deposit@ccdc.cam.ac.uk

To change the displayed columns or filter results you should click on the down arrow of the relevant column heading. To order the results by a particular column click on the column heading you wish to order your results by.

<input type="checkbox"/>	CCDC Nu...	Data block	Deposited On	Deposited By	Refcode	Formula	Embargoed ...	Status	
<input type="checkbox"/>	1416019	data_amit63m	12/01/2017	mplightfoot74@g...	AHUXUY	C14 H13 Cl1 N2 Pd1 S1		Published in the CSD	<input type="button" value="Details"/>
<input type="checkbox"/>	1416026	data_mn467	12/01/2017	mplightfoot74@g...	AHUXOS	C21 H22 O2 S1		Published in the CSD	<input type="button" value="Details"/>
<input type="checkbox"/>	1416025	data_l	12/01/2017	mplightfoot74@g...	AHUXIM	C6 H7 Cl1 N2 O1		Published in the CSD	<input type="button" value="Details"/>
<input type="checkbox"/>	1416016	data_l	10/01/2017	mplightfoot74@g...	AHUXAE	C20 H14 Br1 N3 O3.0.5(C7 H8)		Published in the CSD	<input type="button" value="Details"/>

Over 15,000
registered
users



Enablers of FAIR Crystallography

- **Standard File Formats and Vocabularies**

- Crystallographic Information File (CIF)
- CIF Dictionaries (vocabularies)
- checkCIF Data Validation Service



- **Standard Identifiers – disambiguation and interoperability**

- Digital Object Identifiers – for articles and data
- ORCID iDs – for researchers



- **Trusted searchable data repositories**

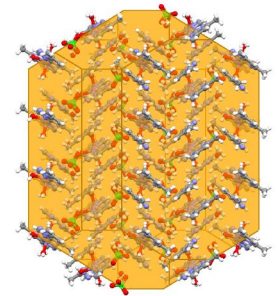
- Cambridge Structural Database
- Inorganic Crystal Structure Database
- Protein Data Bank



Crystal Structure Databases

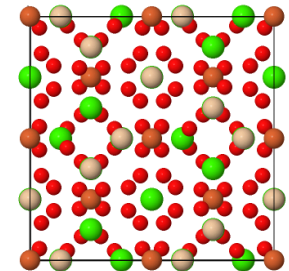
- **Cambridge Structural Database**

- Organic and Metal-organic compounds
- >970,000 structures
- Established in 1965



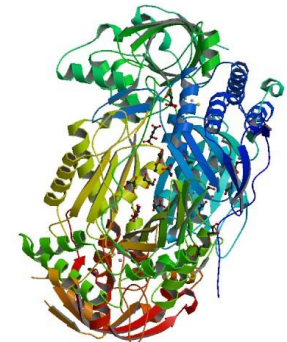
- **Inorganic Crystal Structure Database**

- Inorganic compounds
- >200,000 structures
- Established in 1978



- **Protein Data Bank**

- Biological macromolecules
- >140,000 structures
- Established in 1971



Free, unified deposition and access of crystal structure data

– July 12, 2018

<https://www.ccdc.cam.ac.uk/News/List/2018-07-new-joint-services/>

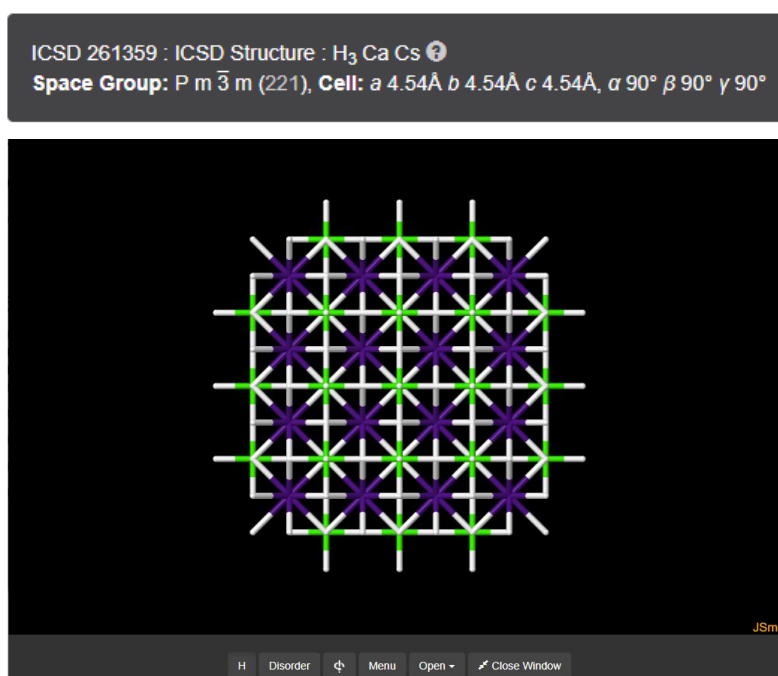
The Cambridge Crystallographic Data Centre (CCDC) and FIZ Karlsruhe – Leibniz Institute for Information Infrastructure (FIZ Karlsruhe) today announced the launch of their joint deposition and access services for crystallographic data across all chemistry. These services will enable researchers to share data through a single deposition portal and explore all chemical structures for free worldwide.



ICSD Entry: 261359

Results	
<input checked="" type="checkbox"/> Database Identifier	Deposition Number
<input checked="" type="checkbox"/> ICSD 261358	1719195
<input checked="" type="checkbox"/> ICSD 261359	1719196
<input checked="" type="checkbox"/> ICSD 261889	1719655
<input checked="" type="checkbox"/> ICSD 261890	1719656
<input checked="" type="checkbox"/> ICSD 422363	1733975
<input checked="" type="checkbox"/> ICSD 422410	1734003
<input checked="" type="checkbox"/> ICSD 422411	1734004

[Download](#)



- New Joint Deposition and Access Portal for Organic and Inorganic crystal structures
- CCDC's infrastructure was adapted to accommodate inorganic structures from FIZ
- Implementation greatly aided by use of a common standard format (CIF)

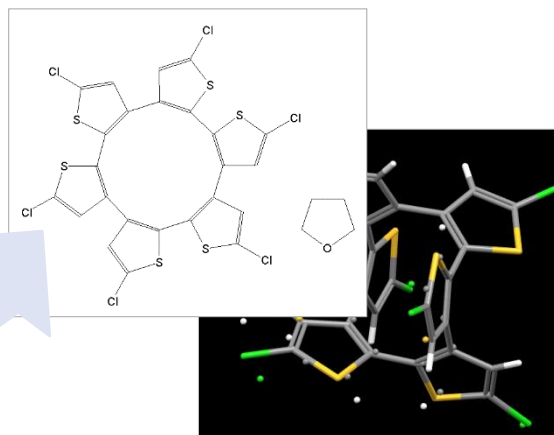


The CSD: Crystallography and Chemistry

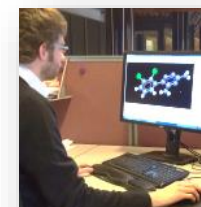
```
loop_  
_atom_site_label  
_atom_site_type_symbol  
_atom_site_fract_x  
_atom_site_fract_y  
_atom_site_fract_z  
_atom_site_U_iso_or_equiv  
_atom_site_adp_type  
_atom_site_occupancy  
_atom_site_symmetry_multiplicity  
_atom_site_calc_flag  
_atom_site_refinement_flags  
_atom_site_disorder_assembly  
_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
```

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



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



A representation of chemistry can be included in a CIF but in deposited files this is rarely found



Chemistry Data Initiatives



Enablers of FAIR Chemistry

- **Technical Enablers**
 - Standard Identifiers (InChI)
 - Open File Formats (Structures)
 - Standard File Formats (Spectra)
 - Terminologies/Vocabularies

- **Social Enablers**
 - Domain Data Activities
 - General Data Initiatives



Linking Crystal Structures to PubChem

PubChem: A database of chemical molecules and their activities against biological assays

InChI Key:

WHGYBXFWUBPSRW-FOUAGVGXSA-N

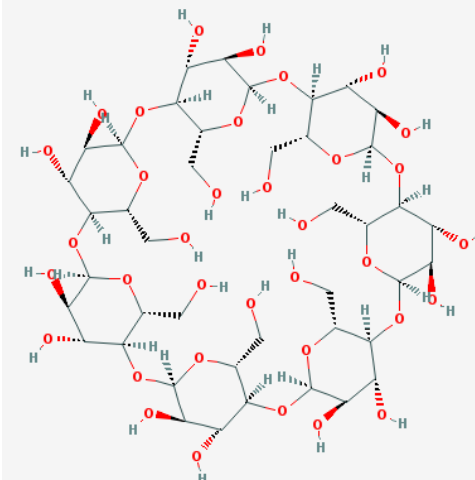
PubChem | OPEN
CHEMISTRY
DATABASE

Compound Summary for CID 444041

PUBCHEM > COMPOUND > BETA-CYCLODEXTRIN

beta-CYCLODEXTRIN

 Vendors  Pharmacology  Literature  Patents  Bioactivities



4.3 Crystal Structures

Crystal Structures: 1 of 1

CCDC Number	762697
Crystal Structure Data	DOI:10.5517/cctln45
Associated Article	DOI:10.1039/C3CE26414A

▶ from The Cambridge Structural Database



Linking Crystal Structures to PubChem

PubChem: A database of chemical molecules and their activities against biological assays



PubChem | OPEN CHEMISTRY DATABASE

Compound Summary for CID 444041

PUBCHEM > COMPOUND > BETA-CYCLODEXTRIN

beta-CYCLODEXTRIN

Vendors Pharmacology Literature Patents Bioactivities

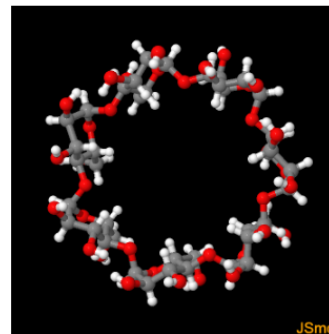
4.3 Crystal Structures

Crystal Structures: 1 of 1

CCDC Number	762697
Crystal Structure Data	DOI:10.5517/cctln45
Associated Article	DOI:10.1039/C3CE26414A

WEWTOJ : 5,10,15,20,25,30,35-heptakis(hydroxymethyl)-2,4,7,9,12,14,17,19,22,24,27,29,32,34-tetradecaooxactacyclo[31.2.2.2^{3,8}.2^{8,11}.2^{13,16}.2^{18,21}.2^{23,26}.2^{28,31}]nonatetracontane-36,37,38,39,40,41,42,43,44,45,46,47,48,49-tetradecol
Space Group: C2, Cell: a 19.056(5)Å b 24.415(6)Å c 15.698(4)Å, α 90.00° β 109.463(13)° γ 90.00°

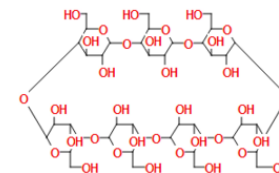
3D viewer



H Disorder Menu Open

Style Labels Packing Measure
Ball and Stick No Labels None None

Chemical diagram



[View group symbols key](#)

Additional CCDC details

CCDC Citation - A.I.Ramos,T.M.Braga,P.Silva,J.A.Fernandes,P.Ribeiro-Claro,M.de F.S.Lopes,F.A.A.Paz,S.S.Braga, CCDC 762697: Experimental Crystal Structure Determination, 2013, DOI: 10.5517/cctln45
Deposited on: 20/1/2010

Associated publications

A.I.Ramos,T.M.Braga,P.Silva,J.A.Fernandes,P.Ribeiro-Claro,M.de F.S.Lopes,F.A.A.Paz,S.S.Braga, *CrystEngComm*, 2013, 15, 2822, DOI: 10.1039/C3CE26414A

from The Cambridge Structural Database



Linking enabled by InChI

PubChem | OPEN CHEMISTRY DATABASE

Bianthrone

Vendors Literature Patents Bioactivities

4.2 Crystal Structures

Crystal Structures: 1 of 16	
CCDC Number	247865
Crystal Structure Data	DOI:10.5517/cc89xn3

from The Cambridge Structural Database

Crystal Structures: 2 of 16	
CCDC Number	299197
Crystal Structure Data	DOI:10.5517/ccb1bj6
Associated Article	DOI:10.1021/jp061205i

from The Cambridge Structural Database

ChemSpider Search and share chemistry

Names and identifiers Properties Searches Spectra **Crystal CIFs** Pharma Links More

Associated Hyperlink: <http://dx.doi.org/10.5517/cc4zfp6>
Comments: Structure CCDC 148418 from the Cambridge Structural Database reported in RSC article <http://dx.doi.org/10.1039/b000825g>
Unit cell: $a=14.4429(5)\text{\AA}$, $b=8.0609(3)\text{\AA}$, $c=24.3908(7)\text{\AA}$, $\alpha=90.00^\circ$, $\beta=99.510(2)^\circ$, $\gamma=90.00^\circ$, $T=123(2)\text{K}$, space group $P21/n$, $Z=4$
Submitted by: [antonywilliams](#)

HM: P21/n
a=14.443Å
b=8.061Å
c=24.391Å
 $\alpha=90.000^\circ$
 $\beta=99.510^\circ$
 $\gamma=90.000^\circ$

IUPAC InChI TRUST

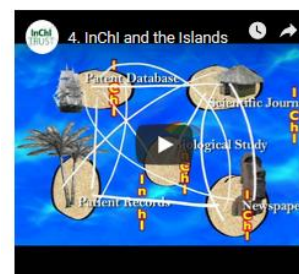
Standard InChI:
InChI=1S/C33H36O6/c1-7-10-37-31-19-25-13-23-17-29(35-5)33(39-12-9-3)21-27(23)15-24-18-30(36-6)32(38-11-8-2)20-26(24)14-22(25)16-28(31)34-4/h7-9,16-21H,1-3,10-15H2,4-6H3

Standard InChIKey:
IZHKSTHBLQRIOW-UHFFFAOYSA-N

The IUPAC International Chemical Identifier

The IUPAC International Chemical Identifier (InChI™) is a non-proprietary identifier for chemical substances that can be used in printed and electronic data sources thus enabling easier linking of diverse data compilations.

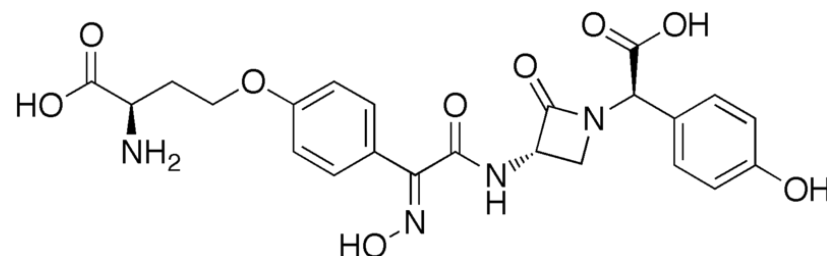
- Initially developed through an IUPAC project from 2000-2004
- Development now overseen by the InChI Trust
- IUPAC still involved in the scientific direction
- Development driven through IUPAC and other Task Groups



<https://www.inchi-trust.org/>



Anatomy of an InChI



Nocardicin A

From Wikipedia, the free encyclopedia

InChI=1S/C23H24N4O9/c24-16(22(31)32)9-10-36-15-7-3-12(4-8-15)18(26-35)20(29)25-17-11-27(21(17)30)19(23(33)34)13-1-5-14(28)6-2-13/h1-8,16-17,19,28,35H,9-11,24H2,(H,25,29)(H,31,32)(H,33,34)/b26-18+/t16-,17+,19-/m1/s1

- **InChI layers**

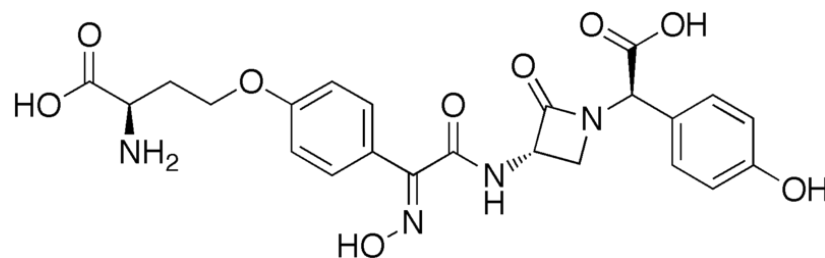
- Main layer
 - formula
 - connections
 - H atoms
- Atom and Bond Stereochemistry
- Isotopic and Fixed Hydrogens (tautomerism)
- Charge



InChIKeys

Nocardicin A

From Wikipedia, the free encyclopedia



InChIKey=CTNZOGJNVIFEBA-TWTPMLPMSA-N

a hashed version of the full standard InChI

- **InChI layers**

- Main

- formula
- connections
- H atoms

- Deprotonation indicator (related to charge)

- Stereo

- Isotopic and Fixed hydrogens

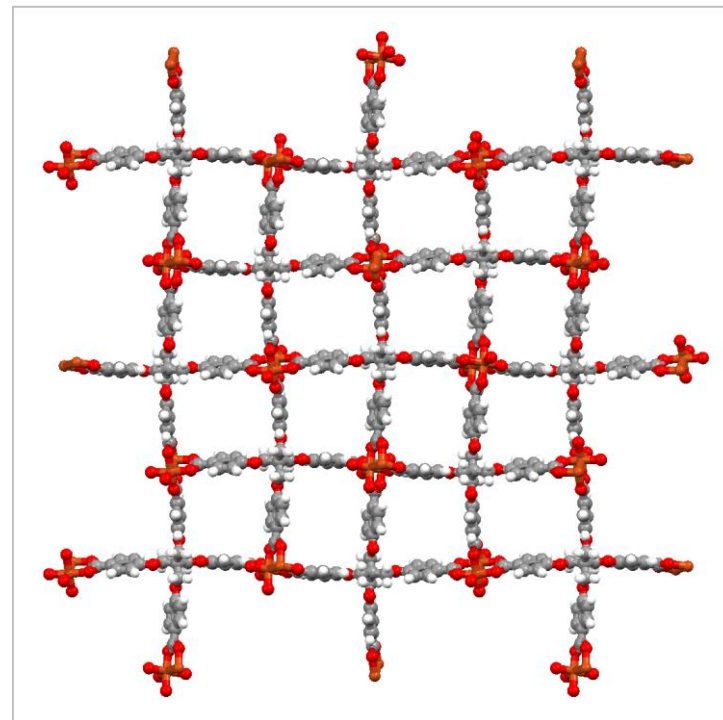
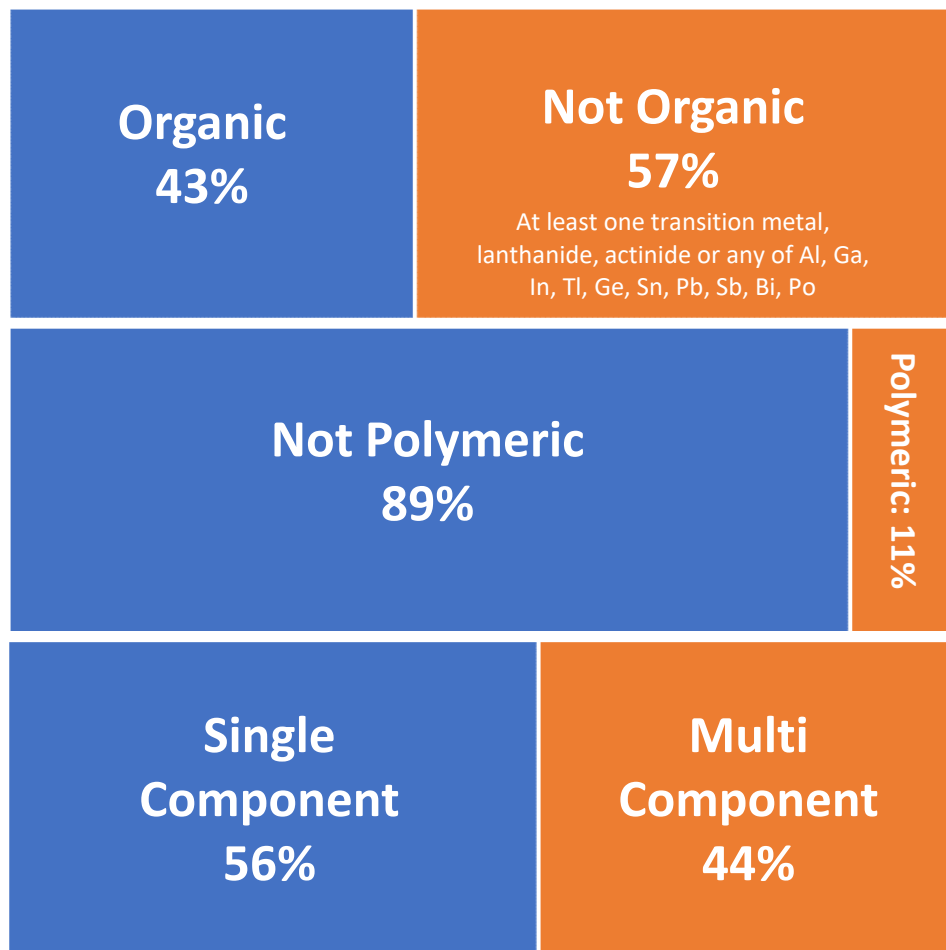
- Standard/non-standard, version number

InChIs can be converted back to structures. InChIKeys cannot.

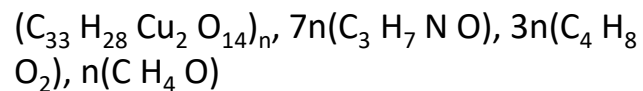
Millions of InChI and InChIKeys have been generated for structures in e.g. PubChem, ChEMBL, ChemSpider



What's in the CSD?



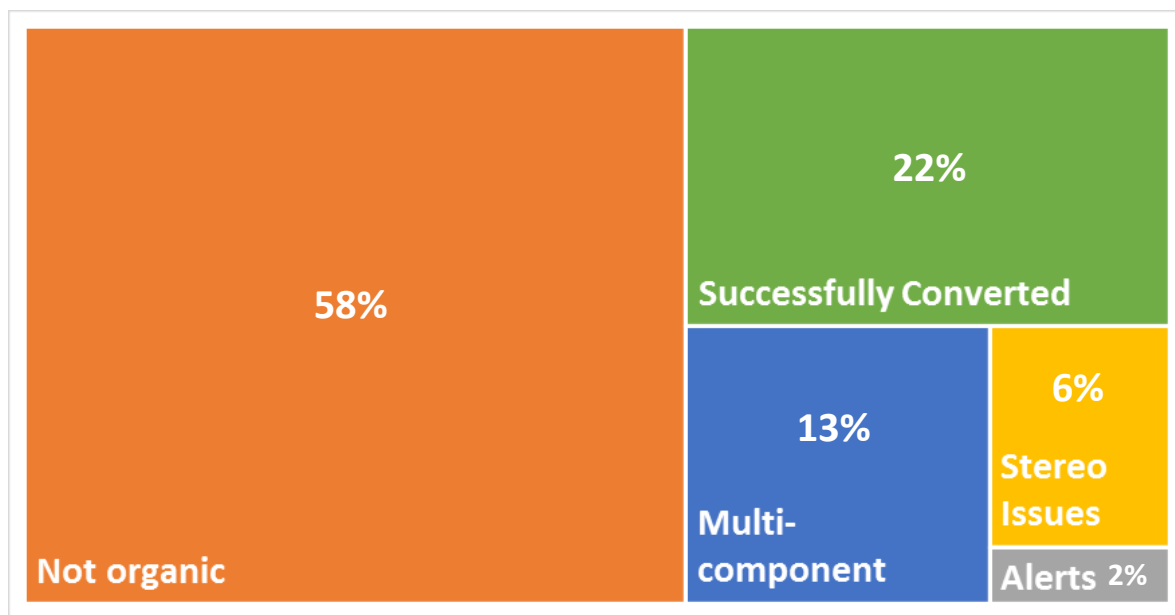
EPOTAF: CCDC 779539 doi:10.5517/ccv55fl





InChI Challenges

Generation of reliable InChIs for chemical substances in the CSD



Based on a subset of 495,751 entries from CSD V5.36

Order of filtering entries out:

- Not organic
- Multi-component
- InChI alerts
- Stereochemistry Issues

- Can confidently generate InChIs for ~22% of CSD entries
- If multi-component entries included then ~35% assuming no other issues
- If based on just organic compounds then 52% (up to 82% including multi-component)

InChI Workshop, Cambridge UK, 4-5 February 2019

Organometallics	Variability/uncertainty, e.g.
Mixtures	- Tautomers
Reactions	- Stereochemistry
QR Codes	- Positional Isomers
Educational Resources	- etc.



InChI Chemical Data Standard: Identifiers and Extensions

Professor Jonathan Goodman: InChI Champion

Wed 13 February 2019, 14:30 - 15:30, Room U203 at the Department of Chemistry. Book at <https://www.training.cam.ac.uk/chem/event/2855805> or email cmc32@cam.ac.uk.

InChI, InChIKeys, Reactions, Mixtures, QR Codes and more...



Chemistry Structure File Formats

- InChIs are generated using a standard InChI Generator supplied by the InChI Trust
- To reliably generate InChIs a **reliable digital representation of a chemical structure** must be supplied as input

- There are **many different file formats** that aim to electronically represent a chemical structure
- Some are ubiquitously used and can be considered ***de facto* standards** – e.g. MOL/SDF, SMILES
- There are ways of **inter-converting** between different file formats

Open Babel: a chemical toolbox designed to speak the many languages of chemical data



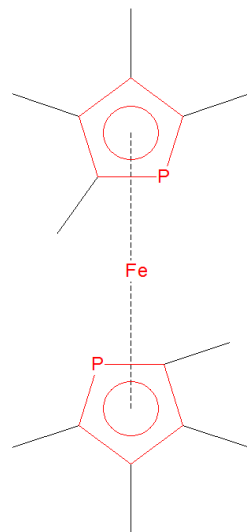
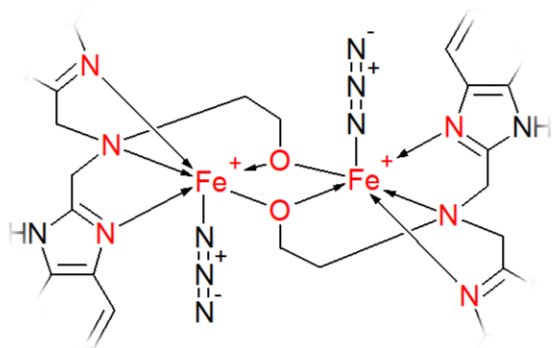
- Common cheminformatics formats
 - Canonical SMILES format (can)
 - Chemical Markup Language (cml, mrv)
 - InChI format (inchi)
 - MDL MOL format (mol, mdl, sdf, sd)
 - Protein Data Bank format (pdb, ent)
 - SMILES format (smi, smiles)
 - Sybyl Mol2 format (ml2, sy2, mol2)
- Other cheminformatics formats
 - Accelrys/MSI Biosym/Insight II CAR format (arc, car)
 - Accelrys/MSI Cerius II MSI format (msi)
 - Accelrys/MSI Quanta CSR format (csr)
 - MCDL format (mcdl)
 - MSI BGF format (bgf)
 - PubChem format (pc)
- Computational chemistry formats
 - ADF cartesian input format (adf)
 - ADF output format (adfout)
 - CAChe MolStruct format (cache, cac)
 - CASTEP format (castep)
 - Cacao Cartesian format (cacrt)
 - Cacao Int
 - DMol3 co

Reads, writes and converts over 110 chemical file formats

Reliable Input Representations

- How best to reliably represent organometallics?

- dative vs covalent bonds?
- explicit hydrogens/valencies?
- dummy atoms?
- zero-order bonds?



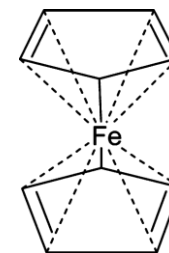
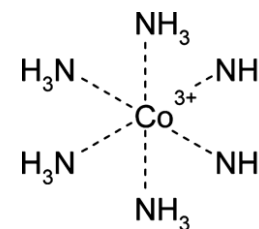
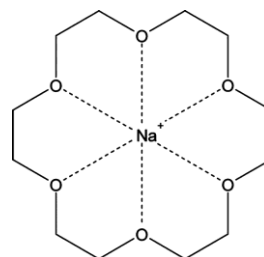
MOL V3000
1 = single
2 = double
3 = triple
9 = coordination
10 = hydrogen
* excluding query bond types

PubChem SDF PUBCHEM_NONSTANDARD BOND
1 Single Bond
2 Double Bond
3 Triple Bond
4 Quadruple Bond
5 Dative Bond
6 Complex Bond
7 Ionic Bond

ACD/Labs MOL V2000 Extensions								
M	ZZF	3	1	41	2	42	3	43
M	ZZH	1	5	2	3	4	5	6
M	ZZH	2	5	7	8	9	10	11
M	ZZH	3	5	12	14	15	16	17
M	ZZE	2	42	18	43	18		

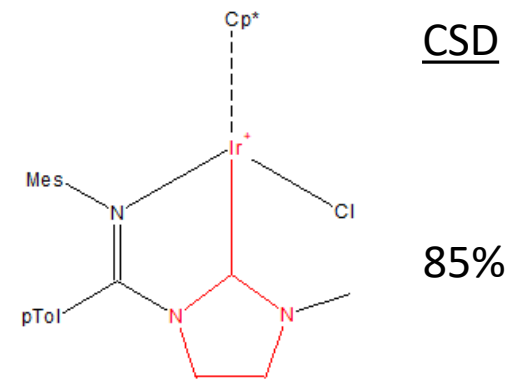
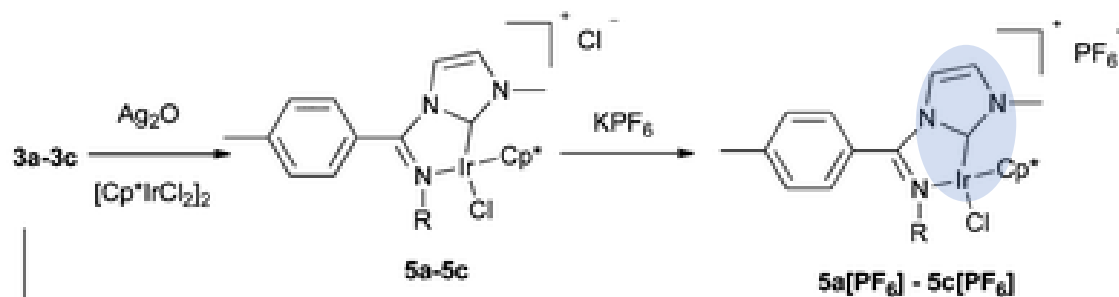
Accurate Specification of Molecular Structures: The Case for Zero-Order Bonds and Explicit Hydrogen Counting. Alex M. Clark.

J. Chem. Inf. Model., 2011, 51 (12), 3149. doi:10.1021/ci200488k



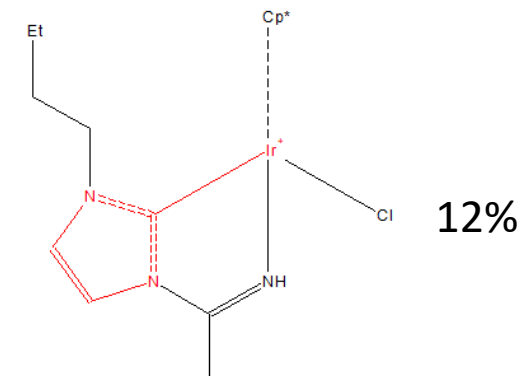
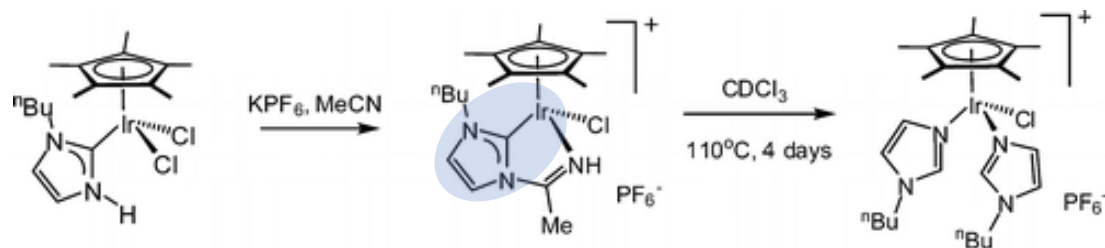


Consistent Structure Representation



ECIWUK CCDC:872879 [10.5517/ccy99dx](https://doi.org/10.5517/ccy99dx)

Dalton Trans., 2012,41, 14557-14567, doi:10.1039/C2DT31989F



LIMXAH CCDC:664254 [10.5517/ccq96kr](https://doi.org/10.5517/ccq96kr)

Organometallics 2007, 26, 18, 4684-4687, doi:10.1021/om700498w



Publishing Chemical Structures

Fragment-Based Approach to Targeting Inosine-5'-monophosphate Dehydrogenase (IMPDH) from *Mycobacterium tuberculosis*

Article Options

Ana Tenen†, Angela Bonitto†, Vinayak Singh§, Mohamed Sabbah†, Anthony C. Cooney†, Valeria

Supporting Information

Fragment-Based Approach to Targeting Inosine-5'-monophosphate Dehydrogenase

	A	B	C
1	Compound_ID	SMILES	Mth IMPDH IC50 (uM)
29	28	<chem>O=C(O)CSC1=NC(C2=CC=CC=C2)=CN1</chem>	–
30	29	https://pubs.acs.org/page/jmcmr/submission/jmcmr_mfstrings.html	
31	30		
32	31		
33	32		
34	33		
35	34		
36	35		
37	36		
38	37		
39	38		
40	39		
41	40		
42	41		
43	42		
44	43		
45	44		
46	45	<chem>O=C(C)NC1=NC(C2=CC=C(Br)C=C2)=CN1</chem>	–
47	46	<chem>O=C(C)NC1=NC(C2=CC=CC=C2)=CN1</chem>	–

Instructions for Authors

1. Use your existing chemical drawing programs (e.g., ChemDraw, ACD ChemSketch, Marvin Sketch) to generate a computer-readable SMILES formula for each compound presented in your article.
2. Paste these formulas into the spreadsheet template, along with basic information about each compound. This spreadsheet will provide a machine-readable version of the key data presented in the article's tables.



Publishing Chemical Spectra

Fragment-Based Approach to Targeting Inosine-5'-monophosphate Dehydrogenase (IMPDH) from *Mycobacterium tuberculosis*

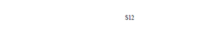
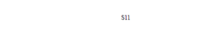
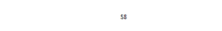
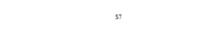
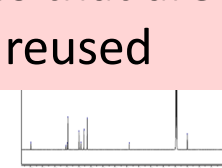
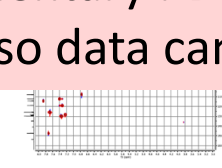
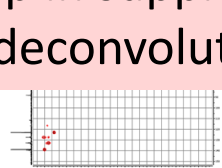
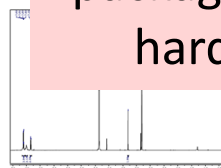
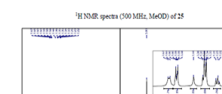
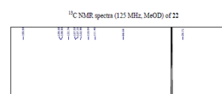
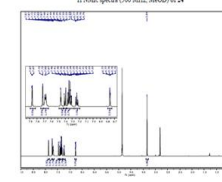
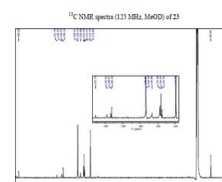
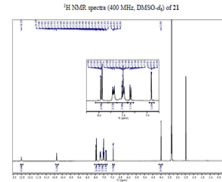
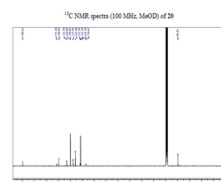
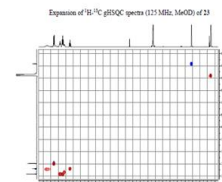
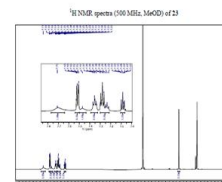
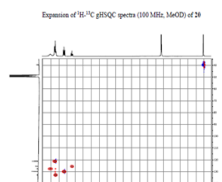
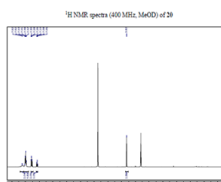
Supporting Information

Fragment-Based Approach

	A	
1	Compound_ID	SMILES
29	28	O=C(O)
30	29	
31	30	
32	31	
33	32	
34	33	
35	34	
36	35	
37	36	
38	37	
39	38	
40	39	
41	40	
42	41	
43	42	
44	43	
45	44	
46	45	O=C(C)
47	46	O=C(C)

Instructions

1. Use your existing data to generate a spreadsheet
2. Paste these spectra into the spreadsheet



Structured data represented as static images and packaged up in supplementary PDF files that are hard to deconvolute so data can be reused

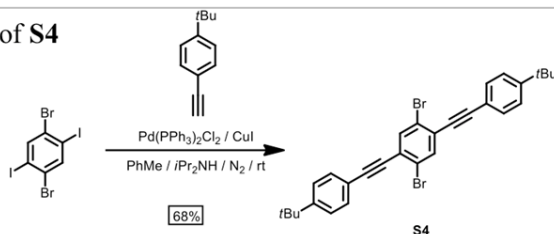
this



Chemistry Data Publication Workflow (simplified)

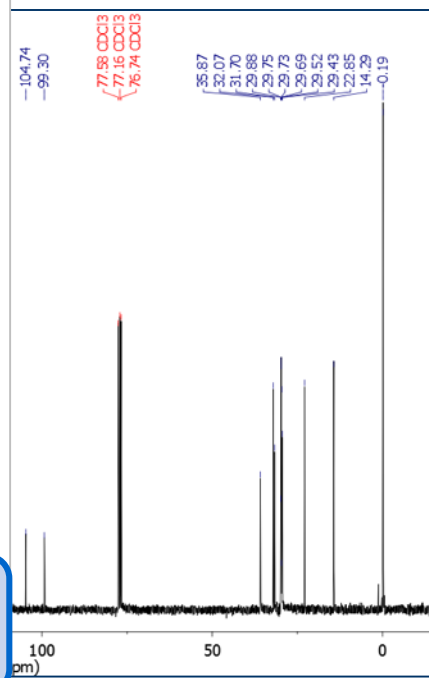


Scheme S6. Synthesis of S4



Synthesis of S4: Anhydrous PhMe (31 mL) and freshly distilled *i*Pr₂NH (10 mL) were added to a 100 mL flask and sparged with N₂ for 20 min. 1,4-dibromo-2,5-diiodobenzene (3.000 g, 6.15 mmol), 1-(*tert*-butyl)-4-ethynylbenzene (2.433 g, 15.38 mmol), Pd(PPh₃)₂Cl₂ (216 mg, 0.308 mmol) and CuI (117 mg, 0.615 mmol) were added to the solution, in sequence. The mixture was stirred at rt for 14 h. The crude reaction mixture was filtered through a pad of SiO₂ gel and washed with additional CH₂Cl₂ (200 mL). Evaporation of the solvent provided the crude product as a solid, which was purified by chromatography (SiO₂, hexanes) to provide S4 (2.310 g, 68% yield) as a white powder. ¹H and ¹³C NMR spectroscopy

were consistent with the report of Hseuh *et al.* S4: ¹H NMR (300 MHz, CDCl₃) δ 7.77 (s, 2H), 7.58 – 7.47 (m, 4H), 7.46 – 7.35 (m, 4H), 1.34 (s, 18H). ¹³C NMR (75 MHz, CDCl₃) δ 152.68, 136.05, 131.70, 126.58, 125.63, 123.81, 119.46, 97.03, 86.49, 77.58, 77.16, 76.74, 35.06, 31.30.

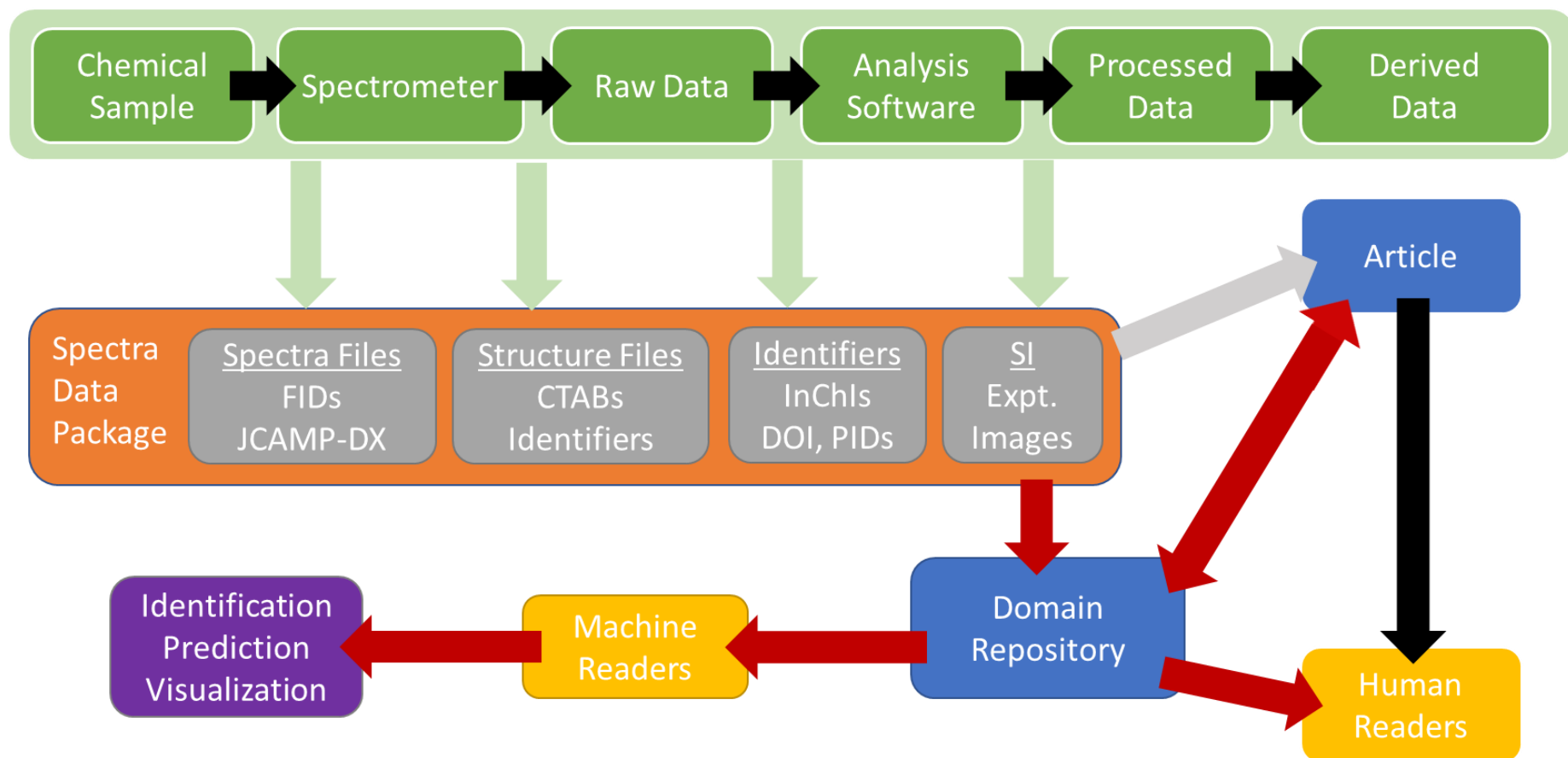


Article

Human Readers



Chemistry Data Publication Workflow (an alternative)





Standards for Spectra: JCAMP-DX

<http://www.jcamp-dx.org/>

- JCAMP-DX is a standard for spectroscopic data developed in the 1980s

JCAMP-DX is a standard file form for exchange of infrared spectra and related chemical and physical information between spectrometer data systems of different manufacture, main-frame time-sharing systems, general purpose lab computers, and personal computers. It is compatible with all media: telephone, magnetic and optical disk, magnetic tape, and even the printed page (via optical reader).

JCAMP-DX: A Standard Form for Exchange of Infrared Spectra in Computer Readable Form
Applied Spectroscopy
(1988)

- It has undergone a number of enhancements to address new instrumentation needs, as well as extensions to new spectroscopic methods.
- Vendors have added custom extensions for their own instruments; new metadata standards, such as ORCIDs and InChIs have been developed.
- Working group of IUPAC Subcommittee on Cheminformatics Data Standards looking at updating the standard

The following FAQs have been asked by members of the Department of Chemistry and answered by members of the Open Data team at the University.

If you have any amendments or further questions you would like to ask please contact the Librarian at the Department of Chemistry, Clair Castle, at library@ch.cam.ac.uk, in the first instance.

The Open Data team can also be contacted at info@data.cam.ac.uk, <http://www.data.cam.ac.uk/>.

FAQs

What would open data for a typical synthetic organic chemistry paper look like?

For a synthetic paper you might include the output files from NMR, UV/Vis, and IR measurements (for example). These should be in a format that others can use, so the data should be in a format that others can use (e.g. CSV for tables, etc.). Images of graphs, especially of NMR experiments, wouldn't meet this criteria). So for example, if you have a table of data (e.g. IR frequencies), you should provide the data in a format that others can use (e.g. CSV).

Lab books form an important record of the experiments performed. At least the detailed methodology for the experiments should be included.

However, if it would be too time consuming and costly to digitise the lab books then you can simply create a meta-data record on the repository so that future users can contact you to physically access your lab books.

What would open data for a typical molecular dynamics based paper look like?

For a computational paper you might include the input and output files from the calculations. Whether you need to include binary output files (which are often produced but hardly ever analysed) is left at your discretion, but if you feel that these files are necessary for the interpretation of the results then they should also be included.

If you have performed a whole suite of experiments, all of which are similar, then it might only be necessary to provide the input files and a couple of example output files. Future researchers can then scrutinize a sample of your output and then re-run all your input files if they wish to do so.

Output files for NMR, UV/Vis, IR etc. should be in formats that others can use – images of graphs, especially of NMR experiments, don't meet this criteria

Convert your files into an open data format

NMR spectroscopy data from TopSpin

These instructions are for converting NMR spectroscopy data from TopSpin to a text file in the internationally accepted open data format JCAMP-DX (<http://www.jcamp-dx.org>).

In TopSpin

File

Save

Save data set in a JCAMP-DX file

OK

Optional: Change name and directory

Leave "Type of archive file = JCAMP DIFF/DUP"

Change "Include these data types =" to "FID+All_PROCNOS"

Leave "JCAMP version = 6.0"

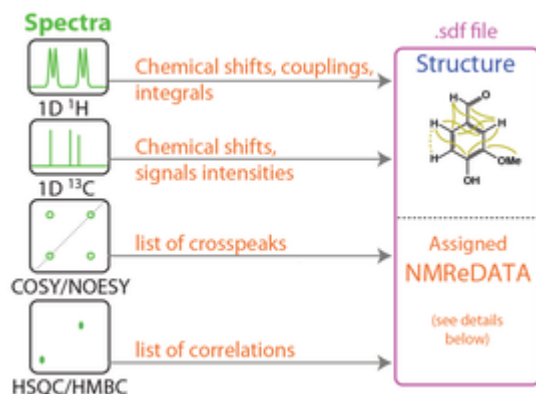
OK

Save as JCAMP-DX – with FIDs

NMReDATA initiative

GENERATE, STORE AND SHARE THE DATA EXTRACTED FROM SET OF NMR SPECTRA ASSOCIATED TO A COMPOUND

The goal of the NMReDATA initiative is to improve the FAIRness and quality of the NMR data available to the community.



Important benefits of the new format

- Improved quality of the NMR data
- Easier inclusion of NMR data in reports and articles
- Simplified referee work
- Compatibility with electronic storage in database
- Easier comparison of dataset
- Improved searchability of NMR data

Mnova Documentation Toolkit – Mpublish

- Researchers with access to Mnova can zip up publication quality images with raw data
- Publisher can digitally sign the zip file so data can be viewed by reviewers and readers without the need for a Mnova licence



<http://resources.mestrelab.com/documentation-toolkit-project-mpublish/>

FAIR Publishing Guidelines for Spectral Data and Chemical Structures in Support of Chemistry and Related Disciplinary Communities

funded by NSF, Orlando FL, March 2019



WORKSHOP GOALS:

1. **Workflow:** develop digital data publishing model across stakeholders
2. **Guidelines:** formulate consistent guidelines for publishing FAIR chemical data for common data types
3. **Value Proposition:** review re-use cases for chemical characterization data
4. **Coalition:** initiate process for ongoing coordination and stakeholder engagement



Publishers • Databases • Repositories • Software Developers • Researchers • Librarians
Standards Organisations • Data Initiatives



Chemistry Publication Guidelines

What data should be published? How should it be validated? Where should it be stored?

- Idea to survey journal requirements for chemistry data initiated at an **IUPAC/RDA workshop**
- Sampling of requirements undertaken by **Vin Scalfani**, University of Alabama Libraries
- Discussion of survey at an **session of the ACS Division of Chemical Information**
- Prompted an **IUPAC taskforce** looking at requirements for a publication of spectra
- Issues relating to file formats discussed at **IUPAC/CODATA workshop** on publishing FAIR data
- Workflow discussions to be advanced at and **NSF Workshop** on Publishing Guidelines for Spectral Data and Chemical Structures





IUPAC Chemical Terminology

- **Blue Book**
 - Nomenclature of Organic Chemistry
- **Red Book**
 - Nomenclature of Inorganic Chemistry
- **White Book**
 - Biochemical Nomenclature
- **Orange Book**
 - Analytical Terminology
- **Purple Book**
 - Compendium of Polymer Terminology and Nomenclature
- **Silver Book**
 - Compendium of Terminology and Nomenclature of Properties Clinical Laboratory Sciences
- **Green Book**
 - Quantities, Units and Symbols in Physical Chemistry





“Digital” Chemical Terminology

<https://goldbook.iupac.org>

- > 7000 terms with authoritative definitions, spanning the whole range of chemistry – with DOIs
- Source documents include *IUPAC Color Books* and recommendations published in *Pure and Applied Chemistry*
- Currently undergoing stabilization and development to provide a foundation for future application

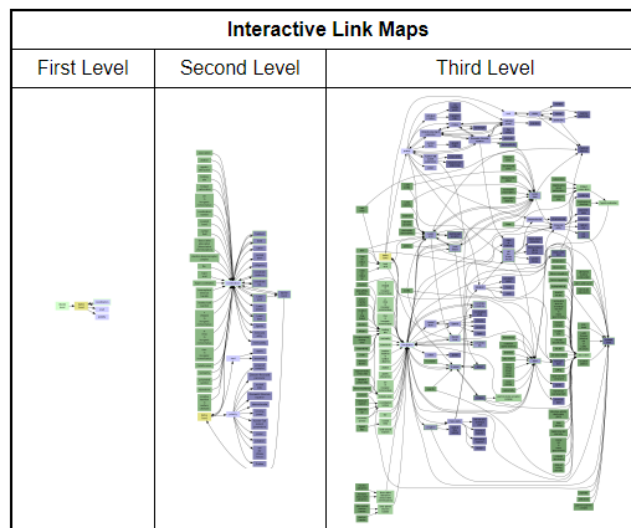
IUPAC GOLD BOOK

dative bond

The coordination bond formed upon interaction between molecular species, one of which serves as a donor and the other as an acceptor of the electron pair to be shared in the complex formed, e.g. the N→B bond in $H_3N \rightarrow BH_3$. In spite of the analogy of dative bonds with covalent bonds, in that both types imply sharing a common electron pair between two vicinal atoms, the former are distinguished by their significant polarity, lesser strength, and greater length. The distinctive feature of dative bonds is that their minimum-energy rupture in the gas phase or in inert solvent follows the heterolytic bond cleavage path.

Source:

PAC, 1999, 71, 1919 (*Glossary of terms used in theoretical organic chemistry*) on page 1933



Cite as:

IUPAC. Compendium of Chemical Terminology, 2nd ed. (the "Gold Book"). Compiled by A. D. McNaught and A. Wilkinson. Blackwell Scientific Publications, Oxford (1997). XML on-line corrected version: <http://goldbook.iupac.org> (2006-) created by M. Nic, J. Jirat, B. Kosata; updates compiled by A. Jenkins. ISBN 0-9678550-9-8. <https://doi.org/10.1351/goldbook>.

Last update: 2014-02-24; version: 2.3.3.

DOI of this term: <https://doi.org/10.1351/goldbook.D01523>.



The Gold Book API Alpha API v0.1 (6/30/17)

While we expect a lot of humans to stop by the Gold Book, its about time that the vocabulary be friendly towards computers and have set up an application programming interface (API) so they may download a bunch of stuff. Here is the overview of the API and we are working on additional documentation. (click the headers below to toggle whats visible.

Terms

Endpoint/Notes

Example(s)

`/terms/index/[scope]/[format]/[download]`

`/terms/index/all` (just "terms" works too)

List of terms in the Gold Book

`/terms/index/C/xml`

[scope]: (all), A-W, XYZ (returns to referring page if no data)

`/terms/index/XYZ/json/download`

[format]: (html), xml, json (rest are ignored)

[download]: (""), download (rest are ignored)

`/terms/view/[identifier]/[format]/[download]`

`/terms/view/A00001`

A term from the Gold Book

`/terms/view/P04409/json`

[identifier]: code, DBid (will expand this)

`/terms/view/ZT07132/xml/download`

[format]: (html), xml, json (rest are ignored)

[download]: (""), download (rest are ignored)

Sources (Click to show)

Ack: S. Chalk



Chemistry Data Initiatives

- **Enablers for FAIR Chemistry data**
 - Standard Identifiers (InChI)
 - Open File Formats (Structures)
 - Standard File Formats (Spectra)
 - Terminologies/Vocabularies
- **Active Communities**
 - InChI Trust
 - IUPAC Committees and Subcommittees
 - ACS Division of Chemical Information (CINF)
 - RDA Chemistry Research Data Interest Group (CRDIG)
 - GO-FAIR Chemistry Implementation Network (ChIN)





Chemistry Data Interest Group: DIGChem

<https://bit.ly/digchem-activity>

DIGChem

Home

DIGChem Events

Publishing Guidelines

Survey

WorkflowTools

JCAMP-DX

NMR/Spectra Repositories

OpenStructures

DataCite Recommendations

Education

Professional Training

Cheminformatics Color Book

Gold Book Website

The Data Interest Group/Chemistry: DIGChem

Coming in 2018:

- [Supporting FAIR Exchange of Chemical Data Through Standards Development](#), Amsterdam, July 16-17, 2018.
- Activities at [ACS/Boston](#), August 19-23, 2018
- [International Data Week](#), Gaborone, Botswana, November 5-8, 2018

[More \(Future/Past\)](#)

The Data Interest Group/Chemistry is an effort to... order to accomplish this vision, CRDIG is analyzing data repositories, evaluating and updating existing... advocating for and educating researchers, librarians... Anyone in the broader chemistry community who is in conjunction with the International Union of Pure and Applied Chemistry (IUPAC), the [International Union of Pure and Applied Chemistry \(IUPAC\) Data Standards \(SCDS\)](#) and the [Research Data Alliance](#). Discussions of the interest group have been held at the International Union of Pure and Applied Chemistry (IUPAC) Information (CINF), at the RSC Chemical Information Society (CINF) General Assembly, at RDA Plenaries, and at the Beijing

Chemistry Research Data IG

IG

Group details

Status: Recognised & Endorsed

Chair (s): Leah McEwen, Stuart Chalk, Ian Bruno, David Martinsen, Richard Kidd



<https://bit.ly/digchem>

Building the social and technical bridges to enable open data sharing



FAIR Data Initiatives



- The **GO FAIR** movement aims to implement the Internet of FAIR Data and Services
- Related to the European Open Science Cloud (EOSC), involving partners outside of the EU

A GO FAIR Chemistry Implementation Network (ChIN) has been recently endorsed

Implementation Networks



Lowering Barriers
to Innovation in
Life Sciences R&D

Implementation of FAIR Data Principles for Pharma and Life Sciences



Kees van Bochove, The Hyve | 07/06/2018 - 05:54

Initiating | Tags: Standards, Knowledge Management



Unmet Needs:

The advent of ML / AI for pharma is very promising, however without a basic amount of metadata and smart annotation of existing data assets, the algorithms cannot make much headway. Several years of IMI knowledge management projects, experiments with data warehouses, data lakes etc. have made it clear that proper semantic annotation of data assets is a hard and resource intensive but very important hurdle to overcome.



Forthcoming Events

- **InChI Workshop: InChI Opportunities**
Department of Chemistry, Cambridge, February 4-5, 2019
- **InChI Chemical Data Standard: Identifiers and Extensions**
Department of Chemistry, Cambridge, February 13, 2019, 14:30 U203
- **FAIR Publishing Guidelines for Spectral Data and Chemical Structures**
Orlando, March 28-29, 2019
- **IUPAC General Assembly and World Congress**
Celebrating 100 Years of IUPAC, Paris, July 5-12, 2019
Special Symposium: Digital Chemistry and the Lab of the Future
- **InChI Symposium**
San Diego, August 23-24, 2019 (Followed by ACS National Meeting)
<https://www.eventbrite.com/e/inchi-symposium-tickets-52810788490>



One Million Crystal Structures: A Wealth of Structural Chemistry Knowledge

Symposium being planned for Fall 2019 ACS Meeting

<https://callforpapers.acs.org/sandiego2019/CINF>



Summary

- **Chemistry Data Initiatives**
 - ❑ Initiatives aimed at supporting FAIR publication of chemistry data
 - ❑ Motivated by guiding principles arising from global data initiatives
 - ❑ Challenges being actively addressed by a range of community groups
- **Cambridge Crystallographic Data Centre**
 - ❑ Sharing crystallographic data and knowledge since 1965
 - ❑ Supporting FAIR publication and access to crystal structures
 - ❑ Adopting and supporting community data standards
- **Cambridge Structural Database**
 - ❑ The world's repository of small molecule crystallographic data
 - ❑ Providing knowledge and insights applicable across chemistry
 - ❑ Available to you at the Department of Chemistry



The Cambridge Crystallographic Data Centre

International Data Repository

Archive of crystal structure data
High quality scientific database

Scientific Software Provider

Search/analysis/visualisation tools
Scientific applications

Collaborative Research Organisation

New methodologies
Fundamental research

Education and Outreach

Conferences, Workshops,
Training, Teaching



@ccdc_cambridge



ccdc.cambridge

<http://www.ccdc.cam.ac.uk/>



*Enriching Chemistry with
Crystallographic Data and Knowledge*

**THANK
YOU!**