



# The Crystallography Open Database – new perspectives

Saulius Gražulis   Andrius Merkys   Antanas Vaitkus   Armel Le Bail  
Daniel Chateigner   Henry Pilliere   Robert T. Downs   Luca Lutterotti  
Peter Moeck   Peter Murray-Rust   Miguel Quirós Olozabal   Werner  
Kaminsky

Denver, SciDataCon2016

Vilnius University Institute of Biotechnology

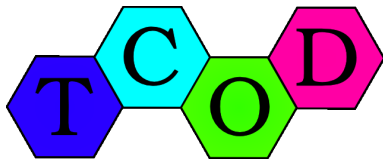


# Open Crystallographic Databases

COD, TCOD, PCOD, MPOD, ...



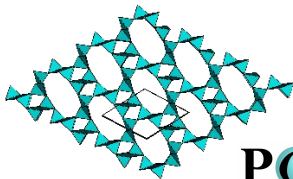
<http://www.crystallography.net/cod>  
> 366 000 entries (ready to grow >  $10^6$ ?)



<http://www.crystallography.net/tcod>  
> 2000 entries (ready to grow to > 350 000?)



<http://mpod.cimav.edu.mx/>  
> 300 entries



<http://www.crystallography.net/pcod>  
>  $10^6$  entries (ready to grow to >  $10^8$ ?)

# The COD project

But what if crystallographers work together to establish a public domain database with all relevant crystallographic data? This would not only overcome the current situation with 'fragmented' databases, it would also prevent for becoming dependent from monopolists.

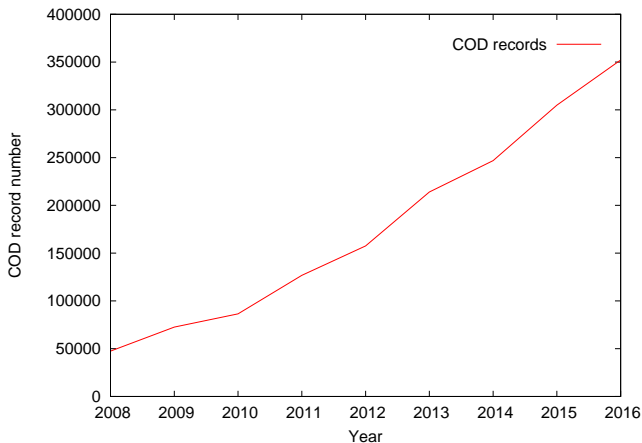
What would be needed?

1. A small team of engaged scientists with some experience in database and software design to coordinate the project.
2. The authors (i.e. the scientific community = YOU) who provides the project with database entries (note, that if you have'nt sold your experimental results exclusively, you are free to distribute the data to such a database, even if they have already been part of a publication - and a lot of good data have never been published).
3. Free software a) for maintaining the database, b) for data evaluation and calculation of derived data (e.g. calculated powder pattern from crystal structures for search-match purposes), c) for browsing and retrieval.

gemstonede (Dr. Michael BERNDT) Fri Feb 14, 2003 1:26 pm

# COD 13 years later

COD increased 7-fold; currently contains over 366000 records (Sept. 2016)



# COD accessibility

COD is a **fully open-access database**. All records are available under public domain designation.

Provided access methods are:

- ▶ Web search
- ▶ URLs constructed from stable identifiers
- ▶ RESTful interfaces
- ▶ Full data download

# COD query examples

Web, REST, SQL

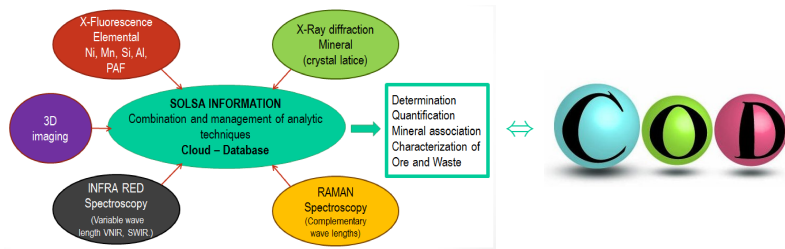
- ▶ Via the WWW interface – go for “search” in:
  - ▶ <http://www.crystallography.net/cod>
  - ▶ <http://www.crystallography.net/tcod>
  - ▶ <http://www.crystallography.net/pcod>
- ▶ Via the **stable** URLs (REST):
  - ▶ <http://www.crystallography.net/cod/2000000.cif>
  - ▶ <http://www.crystallography.net/tcod/10000002.cif>
  - ▶ <http://www.crystallography.net/cod/result?text=perovskite>
- ▶ Via the **views** of the SQL database:
  - ▶ 

```
mysql -u cod_reader cod -h www.crystallography.net \  
-e 'select file, a, b, c, vol, formula  
from data where  
date between "2013-01-01" and  
"2014-12-31" and  
formula regexp " C[0-9]* "  
order by vol desc limit 10'
```

# COD applications

- ▶ SOLSA
  - ▶ <http://www.solsa-mining.eu/>
- ▶ AiiDA [Pizzi et al., 2016]
  - ▶ <http://www.aiida.net/>
- ▶ COSMOS [Sadowski and Baldi, 2013]
  - ▶ <http://cdb.ics.uci.edu/>
- ▶ FPSM [Boullay et al., 2014], MAUD [Boullay et al., 2012]
  - ▶ <http://fpsm.radiographema.com/>
  - ▶ <http://maud.radiographema.eu/>
- ▶ DataWarrior
  - ▶ <http://www.openmolecules.org/datawarrior/>
- ▶ MolView
  - ▶ <http://molview.org/>
- ▶ search-match (Bruker, PANalytical, Rigaku)
- ▶ ... and more!

# SOLSA project and COD



COD will be used in SOLSA for:

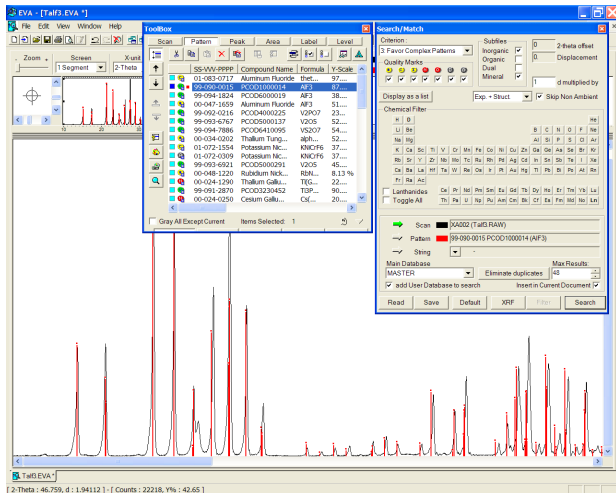
- ▶ mineral identification;
- ▶ subsequent data dissemination.

SOLSA data flow diagram courtesy Monique Le Guen, ERAMET.



# Use of \*COD databases

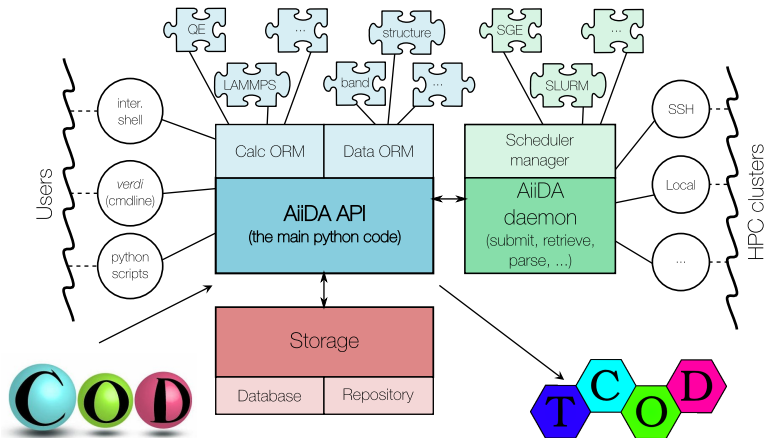
Search-match identification of the materials



A **predicted** phase from PCOD could be identified in experimental data.

Courtesy Armel Le Bail [Le Bail, 2008]

# COD, TCO, and AiiDA link



Courtesy AiiDA developers [Pizzi et al., 2016]

# COD Diffraction Image Store

Uses Tahoe-LAFS (<https://tahoe-lafs.org>) as a back-end:

The screenshot shows a web browser window with the URL <http://192.168.56.101/>. The page title is "Tahoe-LAFS - Welcome". The interface is divided into several sections:

- Grid Status:** Shows the status of the storage grid. It lists "Introducer" (checked) and "Helper" (unchecked). The "Services" section indicates that the "Storage Server" is accepting new shares and is available, while the "Helper" is not running.
- Connected to 2 of 2 known storage servers:** A table showing the status of the storage servers.
- File Management:** Sections for "OPEN TAHOE-URI:", "DOWNLOAD TAHOE-URI:", and "UPLOAD FILE:" with input fields and buttons.

Nickname	Address	Service	Since	Announced	Version	Available
✓ Node2 v0-xu4dy7wvvjraywe346527sh64zrz5q3mk3ahau4pk66iq	(loopback)	storage	12:03:31 11-Sep-2016	12:03:29 11-Sep-2016	allmydata-tahoe1.10.2	2.740B
✓ Node1 v0-1gkuzzmksenqz2wry2hu257y23eweswp05qar1q45ecjrsq	192.168.56.102:43740	storage	12:03:31 11-Sep-2016	12:03:29 11-Sep-2016	allmydata-tahoe1.10.2	2.41GB

Provides:

- ▶ community-backed store ( $\geq 1$  PB)
- ▶ confidentiality through strong encryption
- ▶ extreme hardware loss tolerance

# Interlinked data in COD



## Crystallography Open Database

### COD Home

Home  
What's new?

### Accessing COD Data

Browse  
Search  
Search by structural  
formula

### Add Your Data

Deposit your data  
Manage depositions  
Manage/release  
prepublications

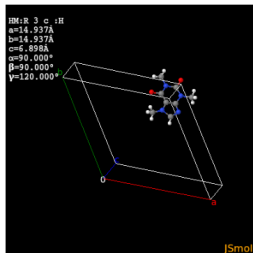
### Documentation

COD Wiki  
Obtaining COD  
Querying COD  
Citing COD  
COD Mirrors  
Advices to donators  
Useful links

### Information card for 2100202

[2100201](#) << [2100202](#) >> [2100203](#)

### Preview



[Display in Jmol](#)

Coordinates

Original IUCr paper

External links

[2100202.cif](#)

[HTML](#)

[ChemSpider](#); [DrugBank](#); [PubChem](#); [Wikipedia](#)

```
select * from wikipedia_x_cod
```

id	ext_id	cod_id	relation_id
1	Ibuprofen	2006278	1
2	Caffeine	2100202	1
3	Serotonin	2019147	1
4	Pristinamycin	1000001	1
5	Cucurbituril	1516465	1
6	Rubrene	1516682	1

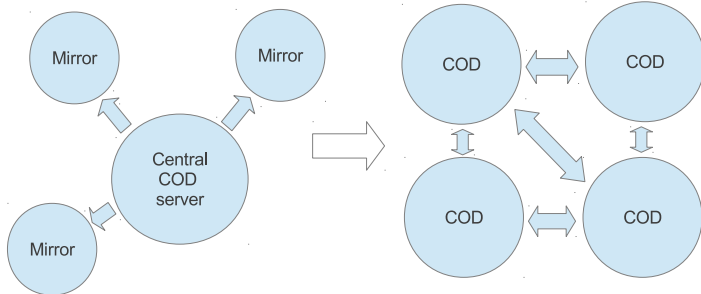
### ▼ Structure parameters

# COD completeness challenge

nr	journal	publisher
45157	Inorganic Chemistry	American Chemical Society
42069	Acta Crystallographica Sect. E	International Union of Crystallography
28775	Dalton transactions (Cambridge ...	Royal Society of Chemistry
26752	Organometallics	American Chemical Society
25493	Journal of the American Chemic ...	American Chemical Society
19824	Acta Crystallographica Sect. C	International Union of Crystallography
19028	Chemical Communications	Royal Society of Chemistry
17858	CrystEngComm	Royal Society of Chemistry
13225	Crystal Growth & Design	American Chemical Society
11083	The Journal of Organic Chemist ...	American Chemical Society
9358	Acta Crystallographica Sect. B	International Union of Crystallography
7910	Organic Letters	American Chemical Society
7516	Dalton Transactions	Royal Society of Chemistry
5751	New Journal of Chemistry	Royal Society of Chemistry
5283	Organic & Biomolecular Chemist ...	Royal Society of Chemistry

# COD durability assurance

- ▶ Best price/performance ratio
- ▶ Capability to build a distributed, equal-peer database



# Acknowledgments

## **VU Institute of Biotechnology**

Virginijus Siksnys  
*(head of the dept.)*

Andrius Merkys  
Antanas Vaitkus

## **QM community**

Björkman  
Torbjörn  
Stefaan Cottenier  
Nicola Marzari  
Giovanni Pizzi  
Lubomir Smrcok  
Linas Vilčiauskas  
Chris Wolverton

## **COD Advisory board**

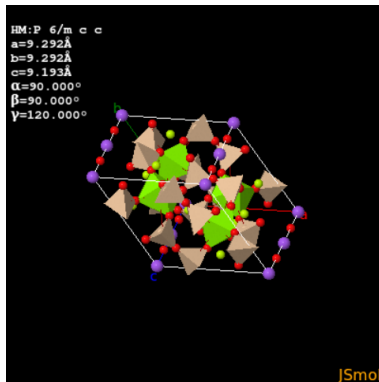
Daniel Chateigner  
Robert T. Downs  
Werner Kaminsky  
Armel Le Bail  
Luca Lutterotti  
Peter Moeck  
Peter Murray-Rust  
Miguel Quirós

Thanks to commercial COD users and supporters – Bruker, PANalytical, Rigaku; thanks to IUCr for support and consultations.

# Thank you!








<http://en.wikipedia.org/wiki/Emerald>



<http://www.crystallography.net/5000095.html>



# References

-  Boullay, P., Lutterotti, L., and Chateigner, D. (2012). Quantitative analysis of electron diffraction ring patterns using the MAUD program.
-  Boullay, P., Lutterotti, L., Chateigner, D., and Sicard, L. (2014). Fast microstructure and phase analyses of nanopowders using combined analysis of transmission electron microscopy scattering patterns. *Acta Crystallographica Section A*, 70:448–456.
-  Le Bail, A. (2008). Frontiers between crystal-structure prediction and determination by powder diffractometry. *Powder Diffraction Suppl.*, pages S5–S12.
-  Pizzi, G., Cepellotti, A., Sabatini, R., Marzari, N., and Kozinsky, B. (2016). AiiDA: automated interactive infrastructure and database for computational science. *Computational Materials Science*, 111:218–230.
-  Sadowski, P. and Baldi, P. (2013). Small-molecule 3d structure prediction using open crystallography data. *Journal of Chemical Information and Modeling*, 53:3127–3130.