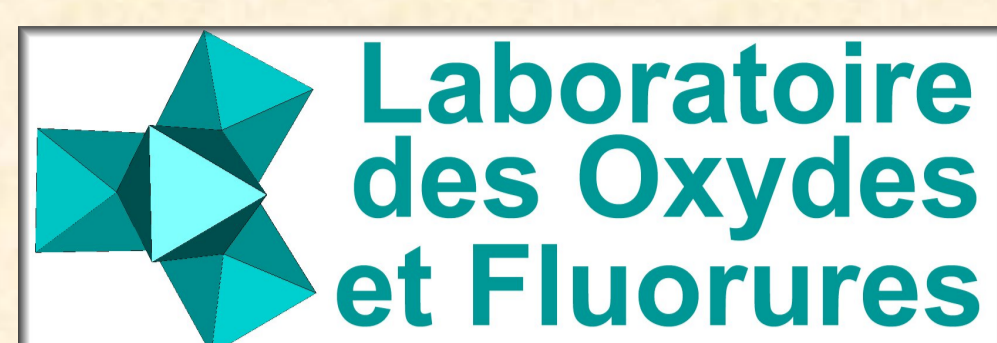


# Software for maintaining and expanding the (Predicted) Crystallography Open Database

Saulius Gražulis<sup>a</sup>, Justas Butkus<sup>a</sup>, Robert Downs<sup>b</sup>, Miguel Quirós Olozabal<sup>c</sup>, Armel Le Bail<sup>d</sup>



<sup>a</sup>Institute of Biotechnology, Graiciuno 8, LT-02241 Vilnius, Lithuania  
<sup>b</sup>Department of Geosciences, University of Arizona, Tucson, Arizona 85721-0077, USA  
<sup>c</sup>Departamento de Química Inorgánica, Facultad de Ciencias, Universidad de Granada, 18071 Granada, Spain  
<sup>d</sup>Université du Maine, Laboratoire des Oxydes et Fluorures, CNRS UMR 6010, Avenue O. Messiaen, 72085 Le Mans Cedex 9, France



Institute of Biotechnology

EU Centre of Excellence

## Introduction

Data are arriving in ever increasing rates. Thus, automated software tools are needed to cope with the growing amount of data, to provide consistent, uniform and accurate information.

The following goals are set by COD team:

- build automated structure deposition tools;
- build a collaboration platform for structure validation and curation;
- ensure data quality – uniformity, integrity, and trustworthiness;
- make scientific data freely accessible to anyone.

## Deposit your data to COD for publication!

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

Crystallography Open Database  
Validation and Deposition Interface

Select CIF file for check:  
home/saulius/ALL.CIF

### About this Validation Interface

This interface allows you to upload, validate and edit CIF files before submitting them for deposition.

### Steps

The process of files deposition, after you have uploaded your data is pretty simple. First step, after files have been uploaded, is validation. Our scripts performs some validation. Results are displayed to you next to your files.

## Automatic data deposition

Crystallography Open Database  
Validation and Deposition Interface

Deposit to COD all valid files

File	Status	Actions
ALL.CIF	valid	<input type="button" value="Edit"/> <input type="button" value="Deposit to COD"/>
File [ALL.CIF] is correct		

You can now check new CIF file.

## CIFParser

A new CIF parser was developed, to meet the following demands:

- identify, locate and clearly show errors in the CIF source (existing parsers were found to be not adequate for this task);
- callable from Perl for ultimate flexibility;
- capable to correct most common CIF errors in a predictable and documented way (via extension a documented and switchable extension of the IUCr CIF grammar);
- easy to extend and check against official CIF grammar for conformance – written in yapp notation, close to a Backus-Naur form.
- free software – distributable under an OSI – compatible license

## cif-tools

available at [svn://www.crystallography.net/cif-tools](http://svn://www.crystallography.net/cif-tools)

A set of tools based on CIFParser was developed and deployed:

- All basic programs a \*x style command programs, developed for Linux;
- agile development – fast deployment, automatic unit testing ensures high quality and predictability
- new tools can be added as needed by anybody;
- integrable into GUIs and Web applications.

Example programs:

cif\_filter – extract uncopywritable data from CIFs, fix errors, compute missing data values;  
cif\_cod\_check – check whether a CIF satisfies COD requirements

## Automatic error detection/correction

Crystallography Open Database  
Validation and Deposition Interface

File	Status	Actions
b407947g-DELIBERATELY-DAMAGED-FOR-TESTING.cif	warnings	<input type="button" value="Edit"/>

The following warnings should be taken into account from file [b407947g-DELIBERATELY-DAMAGED-FOR-TESTING.cif]:

Data block b407947g\_0000001:

- \_journal\_name\_full is undefined
- \_journal\_page\_first is undefined

## Automated CIF download



Remember: atomic coordinates are **not copyrightable**.

The scripts introduce delays so that no noticeable charge will be put in the servers.

COD identifies itself as the downloader.

The bibliographic information is downloaded from Openurl or Pubmed and appended to the CIF's.

## Search by SMILES

```
data_whatever
...
_cell_length_a
...
_atom_site_fract_x
...
_geom_bond_distance
...
```

from CIF to SMILES

```
[Pd]([P](c1cccc1)(C)C(=O)=O)...
```

Use of available free tools: Openbabel (<http://www.openbabel.org>).

Very compact ASCII format

Human easily readable, understandable and editable.

Stores ONLY chemical connectivity, discard any other information.

## SMILES-building helper scripts

Classification of structures: Is it easy to define what is the "molecule" ?

Choosing between openbabel-generated connectivity and author-supplied bonds.

Automatic fixing of some common mistakes.

Removing of chirality markers in non-chiral groups.

Human inspection of the results and fixing if necessary.

## Expanding the PCOD

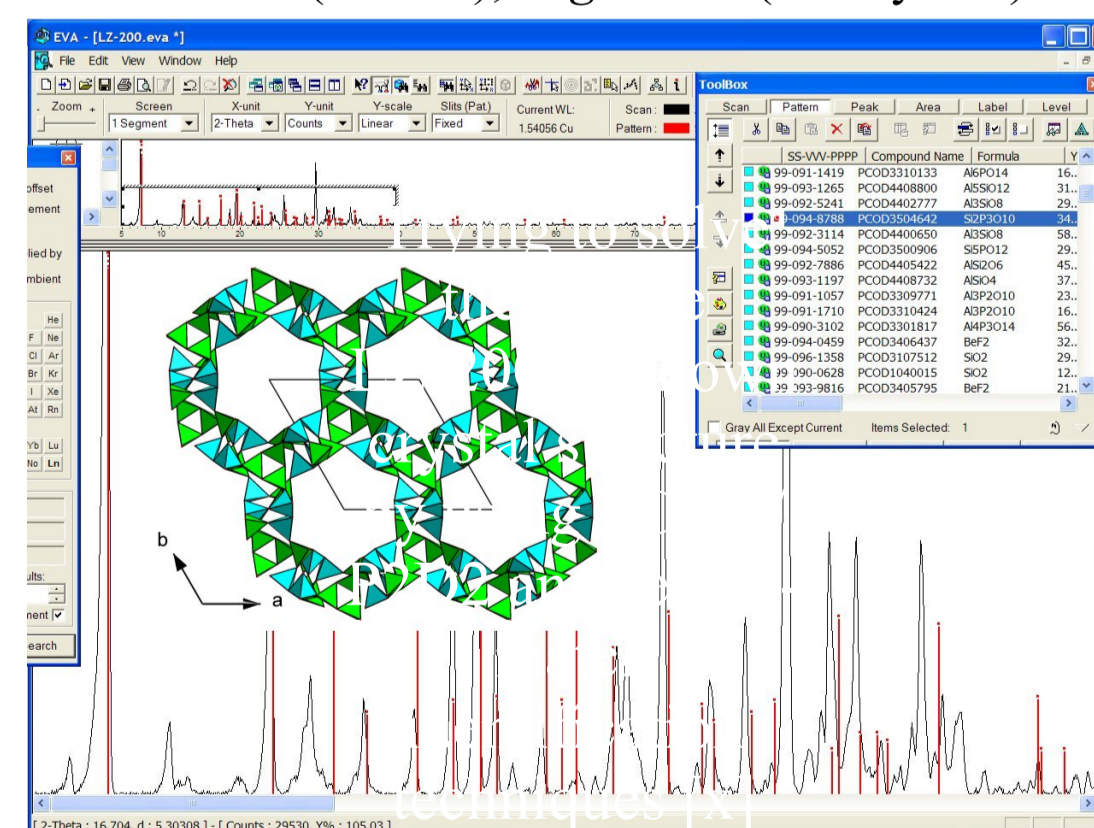
**2010 update** : 898.707 SiO<sub>2</sub> entries were added from ZEFSAII zeolite predictions and the contributions from GRINSP increased to 163.520 (silicates, phosphates, sulfates of Al, Ti, V, Ga, Nb, Zr, or zeolites, fluorides, etc). The PCOD is the first database to attain and offer **more than one million of CIF entries**.

**Software** : a new GRINSP version is now available [3] for parallel computing (for instance using fully the 8 processors of an INTEL core i7).

Other data from other prediction computer programs (CASTEP, CERUS2, CRYSTAL, G42, GULP, USPEX...) are expected, just send them, please.

## PCOD Powder P2D2

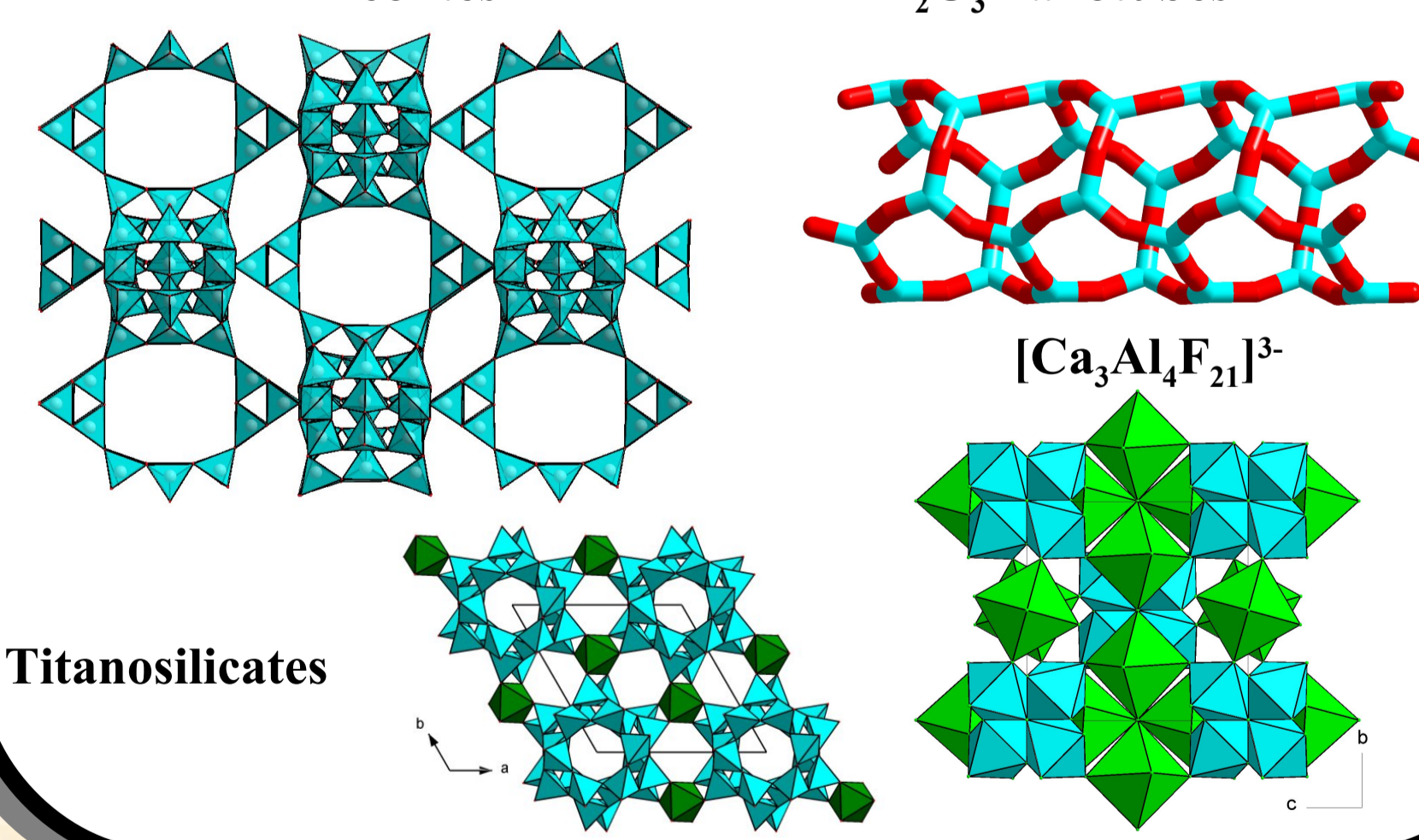
All powder patterns (> 1 million) were calculated and gathered in the P2D2 (Predicted Powder Diffraction Database [4]), they can be used for search-match purposes with EVA (Bruker), Highscore (Panalytical) and more soon.



## VIRTUAL MODELS in PCOD

Zeolites

B<sub>2</sub>O<sub>3</sub> nanotubes



## Acknowledgments

The Vilnius COD development group is financed by the Research Council of Lithuania, contract No. MIP-124/2010

We thank Adriana Daškevič and Andrius Merkys for help with software development and data curation.

Authors thank to all CIF donors, listed on our Web page, and to numerous anonymous volunteers who help to collect data and keep COD running.

The COD Advisory Board thanks commercial supporters for donated hardware and financial support.

## CONCLUSIONS

COD server is technically in position to store and serve all structures that are currently solved. COD deposition procedure ensures syntactic correctness and presence of the most essential CIF data.

COD and PCOD are constantly expanded to meet more needs of crystallographic community

We rely on the help of crystallographic community to add more data and ensure the data correctness!

Use and Add more structures to:  
<http://www.crystallography.net/>

THANK YOU !

## REFERENCES

- [1] <http://www.crystallography.net/>
- [2] S. Gražulis et al., (2009). *J. Appl. Cryst.* 42, 726-729.
- [3] A. Le Bail (2010). *Phys. Chem. Chem. Phys.* 12, 8521-8530.
- [4] A. Le Bail (2008). *Powder Diffr. Suppl.* 23, S5-S12.

COD/PCOD : [cod@crystallography.net](mailto:cod@crystallography.net)