# Open Crystallographic Databases: COD, TCOD and the sisters

Saulius Gražulis

Vilnius, 2023

For the MIF++ seminar, University of Liverpool
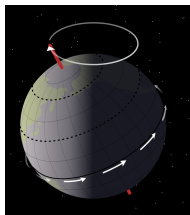
Vilnius University Institute of Biotechnology

# Layout of the talk

1. The value of crystallographic data
2. Crystallographic data(bases): COD, TCOD, PCOD, MPOD, ...
3. Applications of COD and sister databases
4. Mathematical considerations in crystal data processing

# Data importance

**Hipparchus** (c. 190 – c. 120 BCE)

- measured the longitude of Spica and Regulus and other bright stars

- compared his measurements with data from his predecessors, Timocharis and Aristillus, who lived ≈**100** years before him,

- discovered what is now called *the precession of the equinoxes*





By NASA, Public Domain

(Wikipedia, see also articles on Timocharis and Aristyllus)

# Publications are *not* data!

Data need to be extracted (sometimes, manually...) from
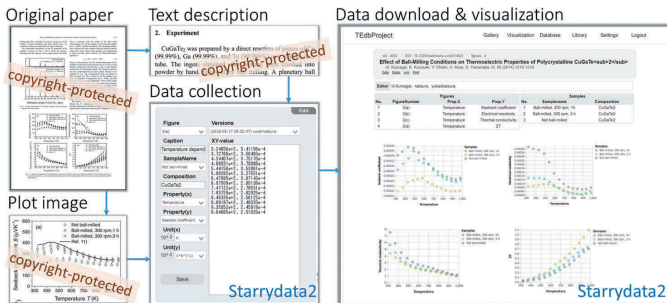publications to make analyses.



**Figure 1.** Concept of plot mining in the *Starrydata2* web system. An example paper [32] and the screenshots of *Starrydata2* web system are presented. Reproduced with permission from Thermoelectrics Society of Japan.

[Katsura et al. (2019)]

# Publications are *not* data!

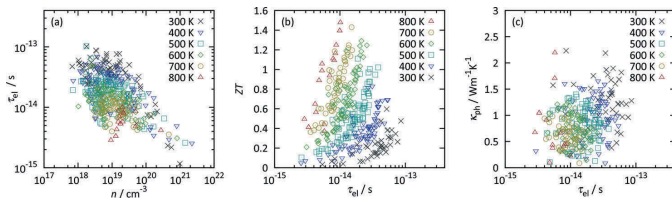But with data, new insights can be drawn from the aggregated publications: https://www.starrydata2.org/



**Figure 6.** Relationship between (a) carrier doping level $n$ and electron relaxation time $\tau_{el}$, (b) $\tau_{el}$ and thermoelectric figure of merit $ZT$, and (c) $\tau_{el}$ and phonon thermal conductivity $\kappa_{ph}$, estimated for 207 experimental samples of $n$-type PbTe.

$$(\tau_{el} \in \left[10^{-15}..10^{-13}\right] \text{ vs. } \tau_{el} = 10^{-14} \text{ s})$$

[Katsura et al. (2019)]

# Crystallographic databases

Open Access:

# Crystallographic databases

Open Access:

- Protein Data Bank;
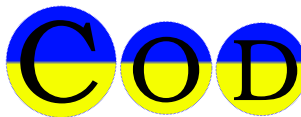
# Crystallographic databases

Open Access:

- Protein Data Bank;

# Crystallographic databases
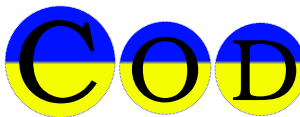
Open Access:

- Protein Data Bank;
- Crystallography Open Database (and its "sisters");

# Crystallographic databases

Open Access:

- Protein Data Bank;
- Crystallography Open Database (and its "sisters");
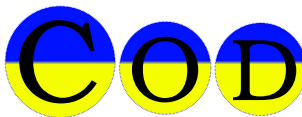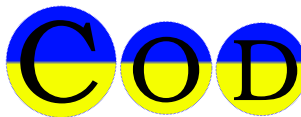- Bilbao Magnetic Structure Database

# Crystallographic databases

Open Access:

- Protein Data Bank;
- Crystallography Open Database (and its "sisters");
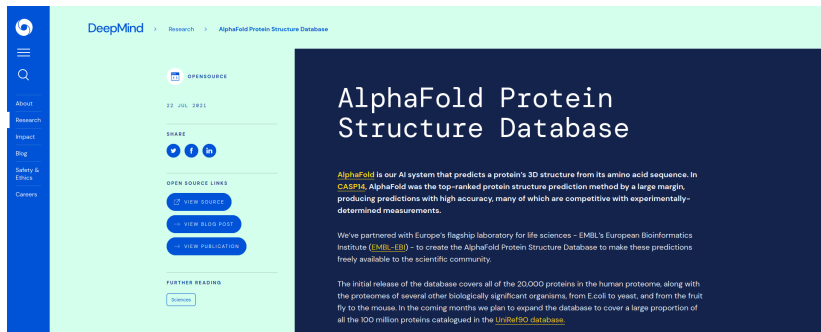- Bilbao Magnetic Structure Database

Proprietary:

- CCDC
- ICSD
- PDF
- Pauling File
- ...

# Crystallographic databases

Open Access:

- Protein Data Bank;
- Crystallography Open Database (and its "sisters");
- Bilbao Magnetic Structure Database

Proprietary:

- CCDC
- ICSD
- PDF
- Pauling File
- ...

About $10^6$ – $10^7$ crystallographic records are available.

# Consequences: AplhaFold

https://deepmind.com/research/open-source/alphafold-protein-structure-database[1]



"Our models are trained on structures extracted from the PDB" [Senior et al. (2020)].

---

[1](accessed 2021-11-23)

# Contents

# 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

# The Crystallography Open Database (COD)

https://www.crystallography.net

Online since 2003 :)

**Crystallography Open Database**

**Open-access collection of crystal structures of organic, inorganic, metal-organic compounds and minerals, excluding biopolymers.**
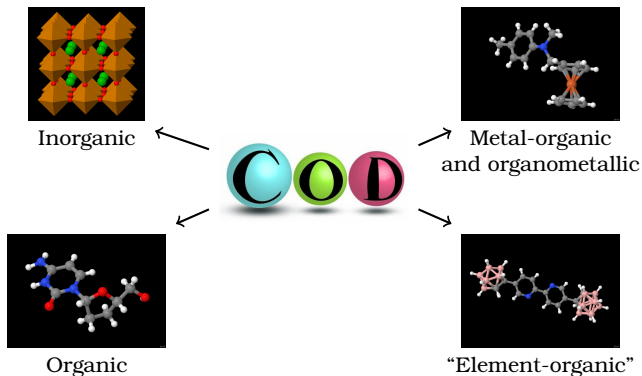
*Including data and software from CrystalEye, developed by Nick Day at the department of Chemistry, the University of Cambridge under supervision of Peter Murray-Rust.*

All data on this site have been placed in the public domain by the contributors.

Currently there are **502408** entries in the COD.

$>$ **500 000** records as of 2023-05-22, available under CC0 License
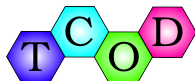
The Crystallography Open Database (COD)
https://www.crystallography.net



Inorganic



Metal-organic
and organometallic



Organic



"Element-organic"

http://www.crystallography.net/cod
> 479 000 entries



http://www.crystallography.net/tcod
> 2900 entries (ready to grow to $> 10^7$?)



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



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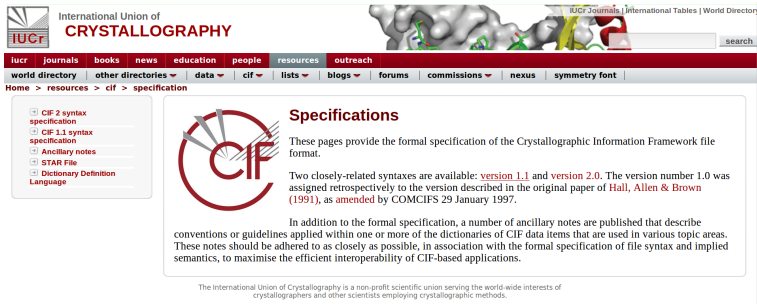


http://solsa.crystallography.net/rod/
> 1100 entries

[Gražulis et al. (2009), Gražulis et al. (2012), Pepponi et al. (2012), Fuentes-Cobas et al. (2017), Mendili et al. (2019)]

# The CIF framework



[Hall et al. (1991)]

The Crystallographic Interchange File/Framework (CIF):

- Provides standard means for data publishing and exchange;
- Is suitable for archiving;
- Is maintained by the IUCr;

# Example of a CIF file

examples/2100858-head.cif:

```
data_2100858
loop_
_publ_author_name
'Buttner, R. H.'
'Maslen, E. N.'
_publ_section_title
;
 Structural parameters and electron difference density in BaTiO~3~
;
_journal_issue                    6
_journal_name_full                'Acta Crystallographica Section B'
_journal_page_first               764
_journal_page_last                769
_journal_volume                   48
_journal_year                     1992
_chemical_compound_source         'synthetic, from a mixture of KF:KMoO4:BaTiO3'
_chemical_formula_sum             'Ba O3 Ti'
_chemical_formula_weight          233.24
_symmetry_cell_setting            tetragonal
_symmetry_space_group_name_Hall   'P 4 -2'
_symmetry_space_group_name_H-M    'P 4 m m'
_cell_angle_alpha                 90.0
_cell_angle_beta                  90.0
_cell_angle_gamma                 90.0
_cell_formula_units_Z             1
_cell_length_a                    3.9998(8)
_cell_length_b                    3.9998(8)
_cell_length_c                    4.0180(8)
```

# CIF atomic coordinates

examples/2100858-coordinates.cif:

```
loop_
_atom_site_label
_atom_site_fract_x
_atom_site_fract_y
_atom_site_fract_z
_atom_site_U_iso_or_equiv
Ba 0.0 0.0 0.0 ?
Ti 0.5 0.5 0.4820(10) ?
O1 0.5 0.5 0.016(5) ?
O2 0.5 0.0 0.515(3) ?
```

# Controlled vocabularies, ontologies

examples/dictionaries/cif-core-example.cif:

```
data_cell_length_
    loop_ _name                     '_cell_length_a'
                                    '_cell_length_b'
                                    '_cell_length_c'
    _category                       cell
    _type                           numb
    _type_conditions                esd
    _enumeration_range              0.0:
    _units                          A
    _units_detail                   'angstroms'
    _definition
;               Unit-cell lengths in angstroms corresponding to the structure
                reported. The values of _refln_index_h, *_k, *_l must
                correspond to the cell defined by these values and _cell_angle_
                values. The values of _diffrn_refln_index_h, *_k, *_l may not
                correspond to these values if a cell transformation took place
                following the measurement of the diffraction intensities. See
                also _diffrn_reflns_transf_matrix_.
;
```

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;
- Better to have no data than wrong data;

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;
- Better to have no data than wrong data;
- Consult original papers or authors themselves if in doubt;

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;
- Better to have no data than wrong data;
- Consult original papers or authors themselves if in doubt;
- Document: record and explain (justify) all changes;

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;
- Better to have no data than wrong data;
- Consult original papers or authors themselves if in doubt;
- Document: record and explain (justify) all changes;
- Keep track of all changes in a version control system;

# COD data management principles

- Strictly stick to IUCr standards (CIF syntax, dictionaries);
- Do not invent data;
- Better to have no data than wrong data;
- Consult original papers or authors themselves if in doubt;
- Document: record and explain (justify) all changes;
- Keep track of all changes in a version control system;
- Keep data provenance (original file names);

# COD data validation

COD data validation policies:

1. Syntactic checks [Merkys et al. (2016)]:
   ```
   $ cifparse 7234818.cif
   ```
2. Semantic validation (against dictionaries)
   [Vaitkus et al. (2021)]:
   ```
   $ cif_validate -D cif_core.dic 7234818.cif
   ```
3. Database-specific checks
   [Gražulis et al. (2009)]:
   ```
   $ cif_cod_check 7234818.cif
   ```

# COD data curation

Data curation in the COD:

```
svn log -r283960 --diff svn://www.crystallography.net/cod/cif/9
```

```
        --- 00/15/9001556.cif (revision 283959)
        +++ 00/15/9001556.cif (revision 283960)
        @@ -68,8 +68,24 @@
        _atom_site_fract_y
        _atom_site_fract_z
        _atom_site_U_iso_or_equiv
        {+_atom_site_type_symbol+}
        {+_atom_site_attached_hydrogens+}
        Fe 0.25000 0.25000 0.25000 0.00490 {+Fe 0+}
        O-H1 0.50000 0.17800 0.30800 0.00100 {+O 1+}
        O-H2 0.19500 0.19000 0.50000 0.00100 {+O 1+}
        O-H3 0.31800 0.50000 0.32300 0.00100 {+O 1+}
        Wat 0.00000 0.50000 0.50000 0.00640 {+O 2+}
        /.../
```

# 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/cod/2000000.html
  - http://www.crystallography.net/cod/result?text=perovskite
- Via the **views** of the SQL database:
  - ```
    mysql -u cod_reader cod -h sql.crystallography.net\
        -e 'select file, a, b, c, vol, formula
            from data where
                year between 2013 and
                              2014 and
                formula regexp " C[0-9]* "
                order by vol desc limit 10'
    ```

# Yes, we OPTIMADE!

http://optimade.org/ [Andersen et al. (2021)]



OPTIMADE
Open Databases Integration
for Materials Design

## http://www.crystallography.net/cod/optimade/v1/structures/

# Contents

# Use of COD and PCOD 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 chemical repertoire

https://molecules.crystallography.net/cod-molecules/cod/2227697.html

SDF file CML file

**Reduced structural formula**



**Reduced canonical SMILES:**

CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C/c1ccc(cc1O)N(CC)CC)(C)C)CC **(x1)** PubChem

## Unique components

| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C /c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)2( /h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

# COD chemical repertoire

https://molecules.crystallography.net/cod-molecules/cod/2227697.html

HM:P -1 #2
a=10.114Å
b=11.400Å
c=13.851Å
α=107.572°
β=110.771°
γ=96.628°

**Reduced structural formula**

A. Vaitkus
ms. ~~in~~
~~preparation~~
under review

SDF file CML file

**Reduced canonical SMILES:**

CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C/c1ccc(cc1O)N(CC)CC)(C)C)CC **(x1)** PubChem

**Unique components**

| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C /c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)2 /h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

# COD chemical repertoire

**Reduced structural formula**

A. Vaitkus
ms. ~~in preparation~~ under review

SDF file CML file

**Reduced canonical SMILES:**

CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C/c1ccc(cc1O)N(CC)CC)(C)C)CC **(x1)** PubChem

**Unique components**

| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C /c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)20(... /h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

https://pubchem.ncbi.nlm.nih.gov/source/849

https://pubchem.ncbi.nlm.nih.gov/substance/164348954

# COD data applications: polymer search

- polymers-in-COD: $\approx 400\,000$ COD records processed
- polymers of different dimensionality (1D, 2D, 3D, 1D-2D and so on) detected, $\approx 93\,000$ polymer records in total.



http://crystallography.net/cod/7224530.html          results of A. Belova

# COD data analysis: polymers

Find interpenetrating chains (crystal nets of covalent bonds):



http://crystallography.net/cod/4103983.html

results of A. Belova, 2019

# COD data analysis: search of knots and links

- Compute knot invariants, such as:
  - linking number;
  - Alexander and/or Conway polynomials;
  - etc. …
- Use the set of invariants to distinguish links and knots.



```
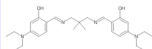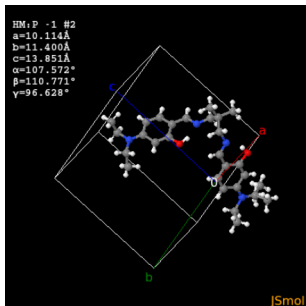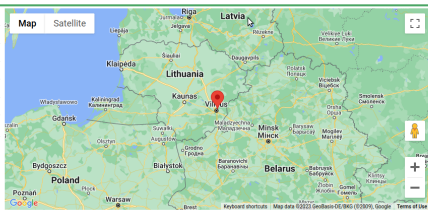# $Id: slides.tex 2298 2023-09-28 12:41:15Z saulius $
#@kw label b1_comp b1_a1 b1_a2 b2_comp b2_a1 b2_a2 sign filename
CROSS X1 C10 a26 a25 C11 a8 a9 -1 hopf-link-integer.cml
CROSS X2 C10 a20 a19 C11 a8 a9 -1 hopf-link-integer.cml
# COMPONENTS: hopf-link-integer.cml 2
# LINKING NUMBER: -1
```

results of A. Belova, 2019

# Contents

# Contents of a crystallographic file

http://www.crystallography.net/cod/2231955.html

# Contents of a crystallographic file

http://www.crystallography.net/cod/2231955.html

# Reconstructing stoichiometric molecular ensemble

1. find a symmetry group $S$ of each molecule;
2. find a symmetry group $H$ the whole molecular ensemble;
3. find (left) coset decomposition of the crystal space group $H$ by $S$, $S \trianglelefteq H$;
4. to each molecule, apply *one* symmetry element from each coset;
5. each choice of symmetry operations from the cosets (transversal) generates a *crystallographically identical* atom set present in the crystal;

# Reconstructing molecules from the COD

http://www.crystallography.net/cod/2231955.html

Usual algorithms:

The new algorithm:

# Reconstructing molecules from the COD

http://www.crystallography.net/cod/2231955.html

Usual algorithms:                          The new algorithm:



[Gražulis et al. (2015)]

A. Vaitkus, cif-perceive-chemistry (formerly cif2molecule) + OpenChemLib

# Predicates on computed data

$IsChiral(M_1) \wedge SymopMapsTo(S, M_1, M_2) \wedge \det(S) = -1$
$\Rightarrow M_1 \; IsEnantiomerOf \; M_2$

# Disorder around a special position

COD 1544968 [Xiang et al. (2016)]

# Disorder around a special position

COD 1544968 [Xiang et al. (2016)]

# Disorder around a special position

COD 1544968 [Xiang et al. (2016)]

# Generating a representative structure

1. find a symmetry group *S* of a special position (Stabiliser);
2. find (left) coset decomposition of the crystal space group *G* by *S*, $S \trianglelefteq G$;
3. take *one* symmetry element from each coset and apply it to the disordered group;
4. each choice of symmetry operations from the cosets (transversal) generates a *distinct* atom set present in the crystal;

Original entry:

http://crystallography.net/cod/4111132.html

# Disorder around a special position in polymers

COD 4111132 [Halper et al. (2006)]



video



video

# Conclusions

- Data publication is as important as papers!
- Aggregated data allows new discoveries...
- ... but for this data need to be properly organised.
- COD, TCOD and the sister databases offer open data in crystallography.
- Mathematical insights are of paramount importance to understand crystal structures.
- Sharing data gives benefits to all.

# Acknowledgements

**VU LSC IBT (KICIS)**

Andrius Merkys
Antanas Vaitkus
Algirdas Grybauskas

**VU LSC IBT (BVTS)**

Daumantas Matulis
Vytautas Petrauskas
Darius Lingė
Marius Gedgaudas

**VU LSC IBT (BNSTS)**

Mindaugas Zaremba
Elena Manakova

**QM community**

Vytautas Žalandauskas
Lukas Razinkovas
Nicola Marzari
Giovanni Pizzi
Lubomir Smrcok
Linas Vilčiauskas
Rickard Armiento

**VU MIF II (FMG)**

Linas Laibinis
Karolis Petrauskas
Haroldas Giedra

**COD Advisory board**

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

**Cheminf community**

Evan Bolton
Paul Thiessen
Thomas Sander

# Thank you!





| | |
|---|---|
| **Coordinates** | [2207377.cif](2207377.cif) |
| **Original IUCr paper** | [HTML](HTML) |

[http://en.wikipedia.org/wiki/Topaz](http://en.wikipedia.org/wiki/Topaz)

[http://www.crystallography.net/2207377.html](http://www.crystallography.net/2207377.html)

Slides available at:
[https://www.crystallography.net/cod/archives/2023/slides/MIF++/slides.pdf](https://www.crystallography.net/cod/archives/2023/slides/MIF++/slides.pdf)

# References I

📄 Andersen CW, Armiento R, Blokhin E, Conduit GJ, Dwaraknath S, Evans ML, et al. (2021) OPTIMADE, an API for exchanging materials data. Scientific Data 8(1):1–10, DOI 10.1038/s41597-021-00974-z, URL https://doi.org/10.1038/s41597-021-00974-z

📄 Fuentes-Cobas LE, Chateigner D, Fuentes-Montero ME, Pepponi G, Grazulis S (2017) The representation of coupling interactions in the Material Properties Open Database (MPOD). Advances in Applied Ceramics 116(8):428–433, DOI 10.1080/17436753.2017.1343782, URL https://doi.org/10.1080/17436753.2017.1343782

📄 Gražulis S, Chateigner D, Downs RT, Yokochi AFT, Quirós M, Lutterotti L, et al. (2009) Crystallography Open Database – an open-access collection of crystal structures. Journal of Applied Crystallography 42:726–729, DOI 10.1107/S0021889809016690, URL http://dx.doi.org/10.1107/S0021889809016690

📄 Gražulis S, Daškevič A, Merkys A, Chateigner D, Lutterotti L, Quirós M, et al. (2012) Crystallography Open Database (COD): an open-access collection of crystal structures and platform for world-wide collaboration. Nucleic Acids Research 40:D420–D427, DOI 10.1093/nar/gkr900, URL http://nar.oxfordjournals.org/content/40/D1/D420.abstract

# References II

Gražulis S, Merkys A, Vaitkus A, Okulič-Kazarinas M (2015) Computing stoichiometric molecular composition from crystal structures. Journal of Applied Crystallography 48:85–91, DOI 10.1107/S1600576714025904, URL http://dx.doi.org/10.1107/S1600576714025904

Hall SR, Allen FH, Brown ID (1991) The crystallographic information file (CIF): a new standard archive file for crystallography. Acta Crystallographica Section A 47:655–685, DOI 10.1107/S010876739101067X, URL http://dx.doi.org/10.1107/S010876739101067X

Halper SR, Do L, Stork JR, Cohen SM (2006) Topological control in heterometallic metal-organic frameworks by anion templating and metalloligand design. Journal of the American Chemical Society 128(47):15,255–15,268, DOI 10.1021/ja0645483, URL https://doi.org/10.1021/ja0645483

Katsura Y, Kumagai M, Kodani T, Kaneshige M, Ando Y, Gunji S, et al. (2019) Data-driven analysis of electron relaxation times in PbTe-type thermoelectric materials. Science and Technology of Advanced Materials 20(1):511–520, DOI 10.1080/14686996.2019.1603885, URL https://doi.org/10.1080/14686996.2019.1603885

Le Bail A (2008) Frontiers between crystal-structure prediction and determination by powder diffractometry. Powder Diffraction Suppl pp S5–S12, DOI 10.1154/1.2903488, URL https://doi.org/10.1154/1.2903488

# References III

📄 Mendili YE, Vaitkus A, Merkys A, Gražulis S, Chateigner D, Mathevet F, et al. (2019) Raman Open Database: first interconnected Raman–X-ray diffraction open-access resource for material identification. Journal of Applied Crystallography 52(3):618–625, DOI 10.1107/s1600576719004229, URL https://doi.org/10.1107/s1600576719004229

📄 Merkys A, Vaitkus A, Butkus J, Okulič-Kazarinas M, Kairys V, Gražulis S (2016) *COD::CIF::Parser*: an error-correcting CIF parser for the Perl language. Journal of Applied Crystallography 49(1):292–301, DOI 10.1107/S1600576715022396, URL http://dx.doi.org/10.1107/S1600576715022396

📄 Pepponi G, Gražulis S, Chateigner D (2012) MPOD: A Material Property Open Database linked to structural information. Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms 284(0):10–14, DOI 10.1016/j.nimb.2011.08.070, URL http://www.sciencedirect.com/science/article/pii/S0168583X11008639, e-MRS 2011 Spring Meeting, Symposium M: X-ray techniques for materials research-from laboratory sources to free electron lasers

# References IV

📄 Senior AW, Evans R, Jumper J, Kirkpatrick J, Sifre L, Green T, et al. (2020) Improved protein structure prediction using potentials from deep learning. Nature 577(7792):706–710, DOI 10.1038/s41586-019-1923-7, URL https://doi.org/10.1038/s41586-019-1923-7, https://doi.org/10.1038/s41586-019-1923-7

📄 Vaitkus A, Merkys A, Gražulis S (2021) Validation of the Crystallography Open Database using the Crystallographic Information Framework. Journal of Applied Crystallography 54(2):1–12, DOI 10.1107/s1600576720016532, URL https://doi.org/10.1107/S1600576720016532

📄 Xiang H, Zhao Q, Tang Z, Xiao J, Xia P, Wang C, et al. (2016) Visible-light-driven, radical-triggered tandem cyclization of o-hydroxyaryl enaminones: Facile access to 3-CF2 /CF3-containing chromones. Organic Letters 19(1):146–149, DOI 10.1021/acs.orglett.6b03441, URL https://doi.org/10.1021/acs.orglett.6b03441