# Sharing scientific data: the crystallography experience

### Saulius Gražulis

## Vilnius, 2023
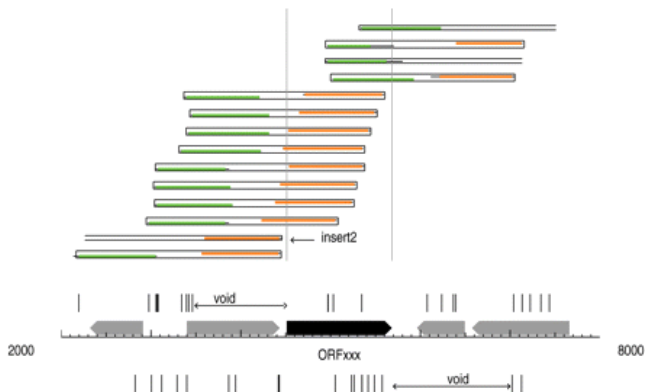
**Vilnius University Institute of Biotechnology**

# Layout of the talk

1. Scientific data: volumes and uses
2. Crystallographic data(bases)
3. Data organisation principles

# Discoveries in raw data

Zheng from the team of Roberts (NEB) use raw sequencing read data to discover *active* restriction-modification systems: [Zheng et al., 2008]:

# Publications are *not* data!
## Starrydata2

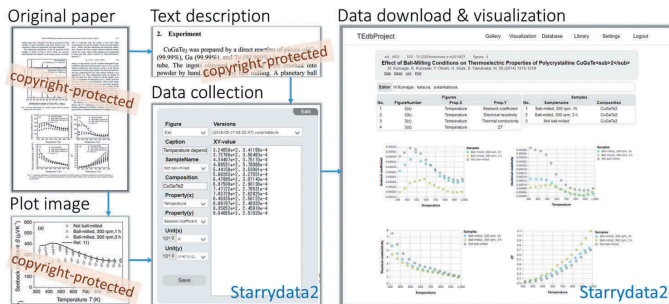Data need to be extracted (sometimes, manually...) from publications to make analyses.



**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/
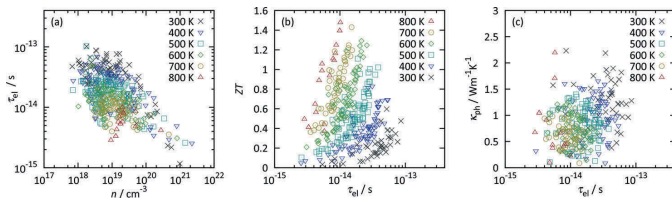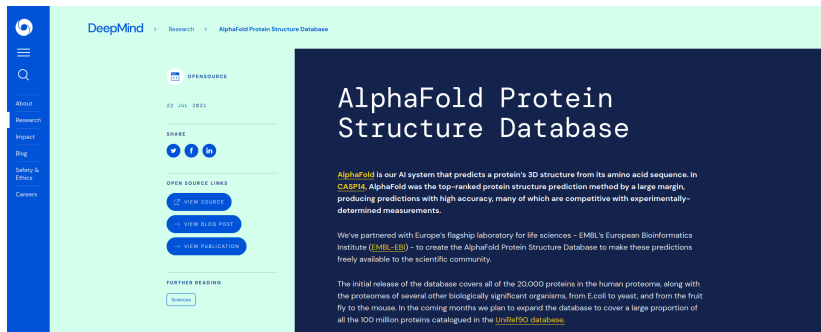


**Figure 6.** Relationship between (a) carrier doping level $n$ and electron relaxation time $\tau_{el}$, (b) $\tau_{el}$ and thermoelectric figure of merit $ZT$, and (c) $\tau_{el}$ and phonon thermal conductivity $\kappa_{ph}$, estimated for 207 experimental samples of $n$-type PbTe.

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

[Katsura et al., 2019]

# 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)

# 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
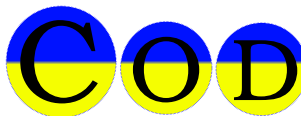- ...

# Crystallographic databases

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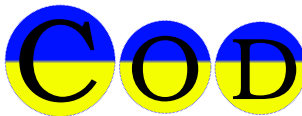
- Protein Data Bank;
- Crystallography Open Database (and its "sisters");
- Bilbao Magnetic Structure Database

Proprietary:

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

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

# The Crystallography Open Database (COD)

https://www.crystallography.net

Online since 2003 :)

**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

**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"

# 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
## CIF (Crystallographic Interchange Framework/Format)

```
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)
```

# 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);

# Three levels of data validation

- Check of file syntax;
- Validation against dictionaries;
- Domain-specific checks:
  - internal consistency;
  - coherence with raw data;
  - scientific plausibility;

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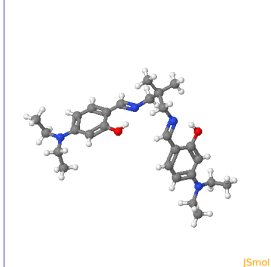
```
svn log -r283960 --diff svn://www.crystallography.net/cod/cif/9

        --- 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+}
        /.../
```
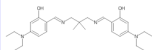
# COD chemical repertoire

jSmol

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**

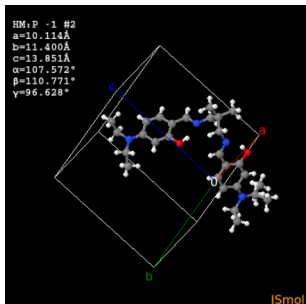| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C<br>/c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)2(...)<br>/h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

See also poster by Merkys et al. (https://bit.ly/3BKZ5vG) in this conference.

# COD chemical repertoire

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



Previous (2227696) Next (2227698) Original COD entry

**Reduced structural formula**

A. Vaitkus
ms. in
preparation

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**

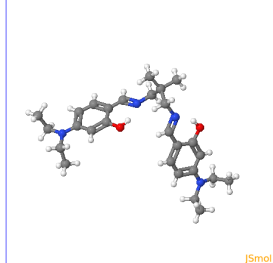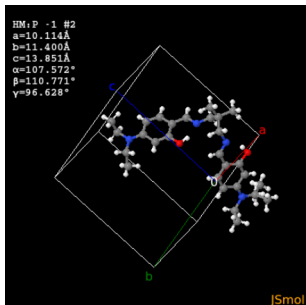| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C<br>/c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)2<br>/h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

See also poster by Merkys et al. (https://bit.ly/3BKZ5vG) in this conference.

# COD chemical repertoire

**Reduced structural formula**

A. Vaitkus
ms. in
preparation

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

Previous (2227696) Next (2227698) Original COD entry

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**

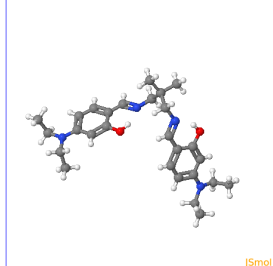| SMILES | InChI |
|---|---|
| CCN(c1ccc(c(c1)O)/C=N/CC(C/N=C /c1ccc(cc1O)N(CC)CC)(C)C)CC | InChI=1S/C27H40N4O2/c1-7-30(8-2)23-13-11-21(25(32)15-23)17-28-19-27(5,6)2( /h11-18,32-33H,7-10,19-20H2,1-6H3/b28-17+,29-18+ |

See also poster by Merkys et al. (https://bit.ly/3BKZ5vG) in this conference.

# COD use cases
## COD and PubChem

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

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

# Conclusions

- Data publication is as important as papers!
- Aggregated data allows new discoveries...
- ... but for this data need to be properly organised.
- Sharing data gives benefits to all.
- **Your contribution is important!**

# 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**

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

**VU MIF II (FMG)**

Linas Laibinis
Karolis Petrauskas

**COD Advisory board**

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

**Cheminf community**

Evan Bolton

# Thank you!





| Coordinates | 2207377.cif |
|---|---|
| Original IUCr paper | HTML |

http://en.wikipedia.org/wiki/Topaz

http://www.crystallography.net/2207377.html

Slides available at:
https://www.crystallography.net/cod/archives/2023/slides/JSMC-2023/slides.pdf

*A path to freedom: GNU → Linux → Ubuntu → MySQL → R → LaTeX → TikZ → Beamer*

Gražulis, S., Chateigner, D., Downs, R. T., Yokochi, A. F. T., Quirós, M., Lutterotti, L., Manakova, E., Butkus, J., Moeck, P., and Le Bail, A. (2009).
Crystallography Open Database – an open-access collection of crystal structures.
*Journal of Applied Crystallography*, 42:726–729.

Hall, S. R., Allen, F. H., and Brown, I. D. (1991).
The crystallographic information file (CIF): a new standard archive file for crystallography.
*Acta Crystallographica Section A*, 47:655–685.

Katsura, Y., Kumagai, M., Kodani, T., Kaneshige, M., Ando, Y., Gunji, S., Imai, Y., Ouchi, H., Tobita, K., Kimura, K., and Tsuda, K. (2019).
Data-driven analysis of electron relaxation times in PbTe-type thermoelectric materials.
*Science and Technology of Advanced Materials*, 20(1):511–520.

Merkys, A., Vaitkus, A., Butkus, J., Okulič-Kazarinas, M., Kairys, V., and Gražulis, S. (2016).
*COD::CIF::Parser*: an error-correcting CIF parser for the Perl language.
*Journal of Applied Crystallography*, 49(1):292–301.

# References II

📄 Senior, A. W., Evans, R., Jumper, J., Kirkpatrick, J., Sifre, L., Green, T., Qin, C., Žídek, A., Nelson, A. W. R., Bridgland, A., Penedones, H., Petersen, S., Simonyan, K., Crossan, S., Kohli, P., Jones, D. T., Silver, D., Kavukcuoglu, K., and Hassabis, D. (2020).
Improved protein structure prediction using potentials from deep learning.
*Nature*, 577(7792):706–710.

📄 Vaitkus, A., Merkys, A., and Gražulis, S. (2021).
Validation of the Crystallography Open Database using the Crystallographic Information Framework.
*Journal of Applied Crystallography*, 54(2):1–12.

📄 Zheng, Y., Posfai, J., Morgan, R. D., Vincze, T., and Roberts, R. J. (2008).
Using shotgun sequence data to find active restriction enzyme genes.
*Nucleic Acids Research*, 37(1):e1–e1.