# Open linked databases in the mining industry

Saulius Gražulis

Kaunas, OpenCon 2017
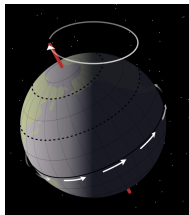
**Vilnius University Institute of Biotechnology**

# 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*





(Wikipedia, see also articles on Timocharis and Aristyllus)

By NASA, Public Domain

# Data and AI systems for geology

[Hart and Duda, 1977]

October 20, 1977

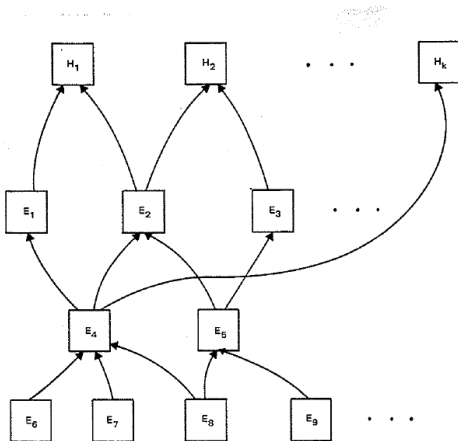PROSPECTOR -- A Computer-Based Consultation
System for Mineral Exploration

by

P. E. Hart and R. O. Duda

Artificial Intelligence Center
SRI International
Menlo Park, California  94025

# The PROSPECTOR network of inference

[Hart and Duda, 1977]
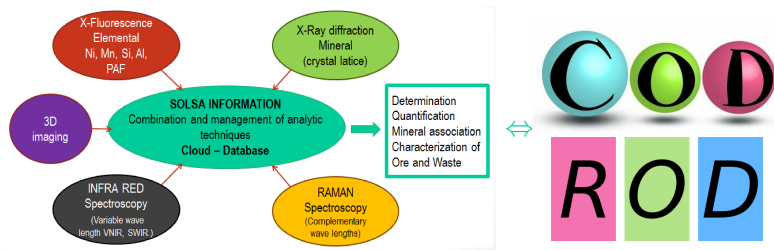
# Data kinds in the SOLSA project



Sonic drilling in Ni laterites

courtesy of Eijkelkamp SonicSampDrill

**Discover SOLSA**

http://solsa-mining.eu/

- ▶ Crystal structures (COD)
- ▶ Raman spectra (ROD)
- ▶ Hyperspectral images (HOD)

# SOLSA project, COD and ROD



COD and other open databases will be used in SOLSA for:

- ▶ mineral identification;
- ▶ subsequent data dissemination.

*SOLSA data flow diagram courtesy Monique Le Guen, ERAMET.*

# Requirements for long-term data archiving and reuse

- Platform independence
  - Text-based formats (ASCII, UTF-8)
- Software independence
- Network-transparency
  - Standard, open protocols (W3C http)
  - Standard, open data carrier formats (JSON, XML, CIF).
  - RESTful servers
- Machine-readable semantics
  - Dictionaries, schemas
- Durability
  - Persistent identifiers
  - Open data principles
  - FAIR principles

# Data exchange in crystallography



[Hall et al., 1991]

The Crystallographic Interchange File/Framework (CIF):

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

# CIF for scientific data

`examples/data/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)
```

# Controlled vocabularies

`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_.
;
```

# Crystallographic data

## The Crystallography Open Database

http://www.crystallography.net/cod

# A COD crystal structure page example
## Sphalerite

http://www.crystallography.net/cod/1525302.html

## Crystallography Open Database

### Information card for entry 1525302

1525301 << **1525302** >> 1525303

**Preview**



HM:F -4 3 m
a=5.427Å
b=5.427Å
c=5.427Å
α=90.000°
β=90.000°
γ=90.000°

Display in Jmol

Coordinates    1525302.cif

Coordinates    1525302.cif

**▼ Structure parameters**

| | |
|---|---|
| Chemical name | (Fe0.2 Mn0.05 Zn0.75) S |
| Formula | Fe0.2 Mn0.05 S Zn0.75 |
| Calculated formula | Fe0.2 Mn0.05 S Zn0.75 |
| Title of publication | Unit-cell edges of natural and synthetic sphalerites |
| Authors of publication | Skinner, B.J. |
| Journal of publication | American Mineralogist |
| Year of publication | 1961 |
| Journal volume | 46 |
| Pages of publication | 1399 - 1411 |
| a | 5.4272 Å |
| b | 5.4272 Å |
| c | 5.4272 Å |
| α | 90° |
| β | 90° |
| γ | 90° |
| Cell volume | 159.855 Å³ |
| Number of distinct elements | 4 |
| Hermann-Mauguin symmetry space group | F -4 3 m |
| Hall symmetry space group | F -4 2 3 |
| Has coordinates | Yes |
| Has disorder | No |
| Has F_obs | No |

### COD Home
Home
What's new?

### Accessing COD Data
Browse
Search
Search by structural formula

### Add Your Data
Deposit your data
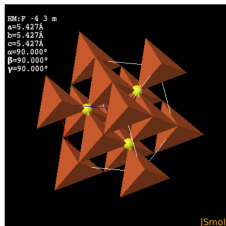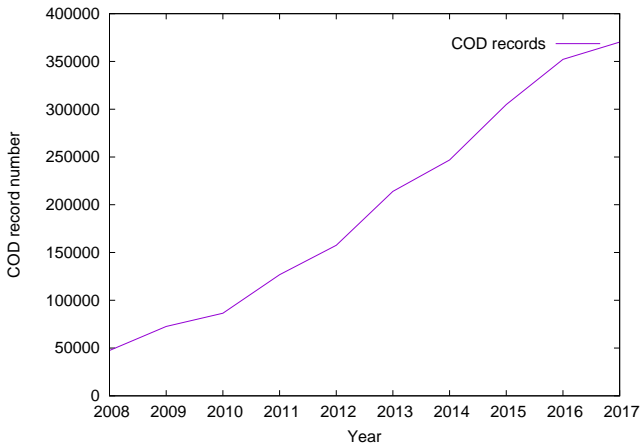Manage depositions
Manage/release prepublications

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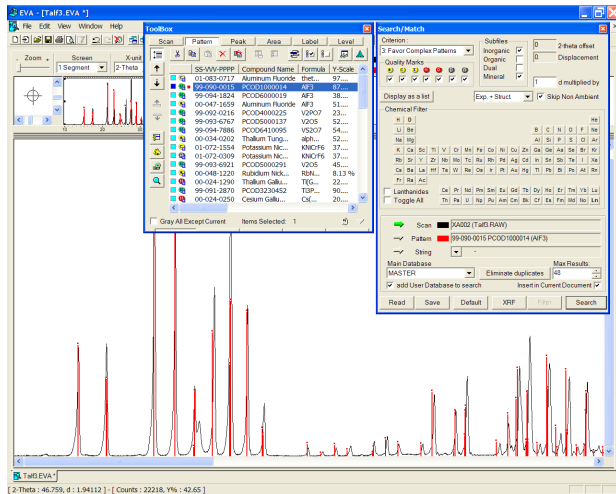COD Mirrors
Advices to donators
Useful links

# COD persistence

COD is on-line for 13 years, increased 7-fold over the last 8 years; currently contains over 385 000 records (October 2017):

# Use of COD and PCOD databases

Search-match identification of the materials



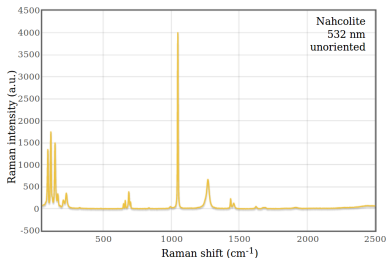A **predicted** phase from PCOD could be identified in experimental data.

Courtesy Armel Le Bail
[Le Bail, 2008]

# Ramano spectroscopy



us-tech.co.za



ROD 3500101

- ▶ the method is very fast
- ▶ requires comprehensive, high quality database

# Raman spectroscopy data

## The Raman Open Database

http://solsa.crystallography.net/rod

*Data records contributed to the ROD by Yassine El Mendili*

# ROD data files

### ROD uses CIF syntax

`examples/data/3500024-head.rod:`

```
#------------------------------------------------------------------------------
#$Date: 2017-10-05 18:15:36 +0300 (Thu, 05 Oct 2017) $
#$Revision: 219 $
#$URL: svn://172.16.1.102/rod/cif/3/50/00/3500024.rod $
#------------------------------------------------------------------------------
#
# This file is available in the Raman Open Database (ROD),
# http://solsa.crystallography.net/rod/
#
# All data on this site have been placed in the public domain by the
# contributors.
#
data_3500024
loop_
_publ_author_name
'El Mendili, Y'
_publ_section_title
;
 SOLSA communication to ROD
;
_journal_name_full           'Personal communication to ROD'
_journal_year                2017
_chemical_compound_source    'commercial powder Prolabo pur'
_chemical_formula_structural 'O2 Ti'
```

# The ROD dictionary

## ROD uses controlled CIF vocabulary

http://solsa.crystallography.net/rod/cif/dictionaries/cif_raman_0.1.1.dic
http://solsa.crystallography.net/rod/cif/dictionaries/cif_rod_0.1.0.dic

examples/dictionaries/raman-example.dic:

```
save__raman_measurement_device.direction_polarization
    _definition.id             '_raman_measurement_device.direction_polarization'
# ... some text omited for brevity ...
    _definition.update          2017-04-10
    _description.text
;
    The direction polarization of the measurement device.
;
# ...
    loop_
    _enumeration_set.state
    _enumeration_set.detail
    unoriented
;
 Unoriented.
;
    Z(XX)Z
;
 Laser polarized parallel to the x axis; analyzer set to pass the x axis
 polarized light.
;
```

*ROD dictionaries coded by Antanas Vaitkus*

# Semantic versioning of the ROD dictionaries

- ROD dictionaries undergo semantic versioning:
  - Bug-fix releases (1.2.x) are compatible backwards and forward;
  - Minor releases (1.x) are backwards compatible;
  - Incompatible changes will be marked by major releases (1.x $\to$ 2.x);

# COD query examples

### Web, REST, SQL

- ▶ Via the WWW interface – go for "search" in:
    - ▶ http://www.crystallography.net/cod
    - ▶ http://solsa.crystallography.net/rod
    - ▶ http://solsa.crystallography.net/hod
- ▶ Via the **stable** URLs (REST):
    - ▶ http://www.crystallography.net/cod/2000000.cif
    - ▶ http://solsa.crystallography.net/rod/3500021.rod
    - ▶ http://solsa.crystallography.net/rod/3500021.html
    - ▶ 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
                year between 2013 and
                              2014 and
                formula regexp " C[0-9]* "
                order by vol desc limit 10'
    ```

# Open Crystallographic Databases

## COD, TCOD, PCOD, MPOD, ROD, HOD ...



http://www.crystallography.net/cod

> 385 000 entries (ready to grow $> 10^6$?)



http://www.crystallography.net/tcod

> 2500 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/

> 120 entries

## HOD

http://solsa.crystallography.net/hod/

TBA...

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

# Hyperspectral image database (HOD)

http://solsa.crystallography.net/hod

A "hybrid" approach necessary due to large size of raster data:

- Metadata and image headers stored in CIF;
- Raster data stored as "raw" binaries;

# HOD record example

examples/hod/1000000-head.cif:

```
data_1000000
loop_
_[local]_description
'ENVI File'
'Created [Wed Jun 08 12:34:07 2016]'
_[local]_wavelength_units        Nanometers
loop_
_hyper_bands.default
220
227
253
_hyper_bands.lines               937
_hyper_bands.number              288
_hyper_bands.samples             384
_hyper_file.byte_order           0
_hyper_file.data_type            4
_hyper_file.type                 ENVI_Standard
_hyper_header.offset             0
_hyper_header_file.contents
;ENVI
description = {
  ENVI File, Created [Wed Jun 08 12:34:07 2016]}
samples = 384
lines   = 937
```



**Test Hyperspectral Open Database**

**Information card for entry 1000000**

4060001 << **1000000** >> 4060060
Search

**Preview**

# SOLSA Large File Store

### Suitable, e.g., for images

Uses Tahoe-LAFS (https://tahoe-lafs.org) as a
back-end [Selimi and Freitag, 2014]:

Tahoe-LAFS architecture



Tahoe-LAFS
storage servers,
direct attached storage

Tahoe-LAFS
storage protocol
over SSL

Tahoe-LAFS gateway

Tahoe-LAFS
storage
client

HTTP(S)
server

Tahoe-LAFS
REST
web-API

over HTTP(S)
or (S)FTP

Tahoe-LAFS client

• web browser
• command-line tool
• Windows virtual drive
• JavaScript frontends
• tahoe backup tool
• duplicity
• (S)FTP client
• FUSE

security perimeter for confidentiality and integrity

Red means that whoever controls that link or that machine can
see and change the contents of your files. You *rely on* that
component for confidentiality and integrity.

Black means that control of that link or that machine does not
give the ability to see or change the contents of your files.
You *do not rely on* that component for confidentiality or
integrity.

Quoted from https://tahoe-lafs.org/trac/tahoe-lafs

# Tahoe LAFS Grid for SOLSA

*Tahoe-LAFS for SOLSA set up by Erikas Raginis*
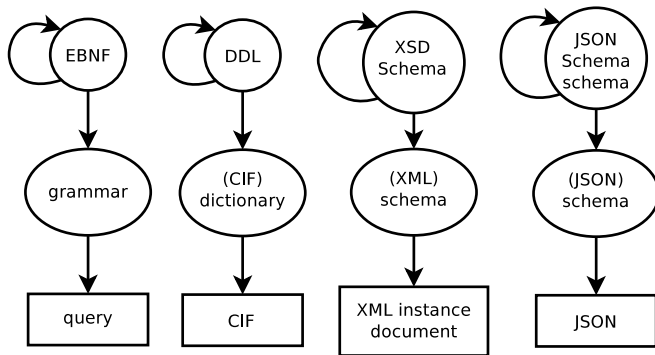
# HOD files on the Tahoe LAFS grid

# HOD (large) data retention policy

A managed data phase-out policy possible:

- Keep data that are:
    - The first of their kind;
    - The best of their kind;
    - The most often used/cited;
    - A small but representative test set (for software);
- Apply lossy compression to older records ($\times 20$ fold possible)
- Discard data for other records, leave just (aggregated) metadata;

# Common pattern of self-describing data definitions

# Acknowledgements

**VU Institute of Biotechnology**

Virginijus Siksnys
( *head of the dept.* )

Andrius Merkys
Antanas Vaitkus
Erikas Raginis

**The SOLSA team**

Monique Le Guen
Beate Orberger
Daniel Chateigner
Henry Pilliere
*and all the team working on the project!*

**COD Advisory Board**

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

# Thank you!



```
HM:P 42/m n m
a=4.594Å
b=4.594Å
c=2.959Å
α=90.000°
β=90.000°
γ=90.000°
```

JSmol

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

*Rob Lavinsky, iRocks.com – CC-BY-SA-3.0*

http://www.crystallography.net/9015662.html

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

# References I

📄 Fielding, R. T. (2000).
*Architectural Styles and the Design of Network-based Software Architectures*.
PhD thesis, University of California, Irvine.

📄 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.

📄 Hart, P. E. and Duda, R. O. (1977).
Prospector – a computer-based consultation system for mineral exploration.
techreport, Artificial Intelligence Center, SRI International, Menlo Park, California 94025.

# References II

📄 Le Bail, A. (2008).
Frontiers between crystal-structure prediction and
determination by powder diffractometry.
*Powder Diffraction Suppl.*, pages S5–S12.

📄 Selimi, M. and Freitag, F. (2014).
Tahoe-lafs distributed storage service in community
network clouds.
*2014 IEEE Fourth International Conference on Big
Data and Cloud Computing.*

RESTful queries [Fielding, 2000]:

- Programming language, transfer protocol **independent**
- GET queries should be null-potent (do not change anything; always provide the same result for the same query);
- POST/PUT queries should be idempotent (the same query executed several times should have the same result as just one query).