# Efficient long-term open-access data archiving in mining industries

<u>Saulius Gražulis</u> & the SOLSA consortium

## Amsterdam, RTM Conference, 2017
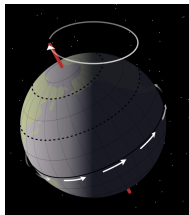
Vilnius University Institute of Biotechnology

# 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

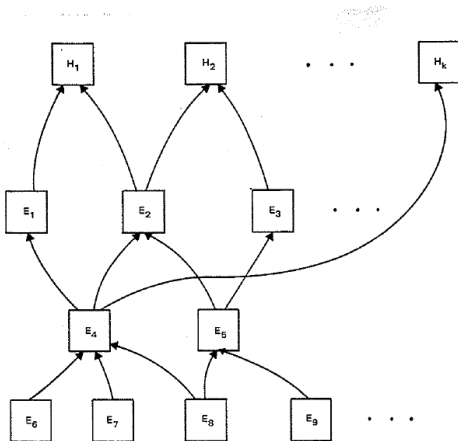PROSPECTOR -- A Computer-Based Consultation
System for Mineral Exploration

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



**Discover SOLSA**

http://solsa-mining.eu/

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

# 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 archiving;
▶ 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)
```

SOLSA

# 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

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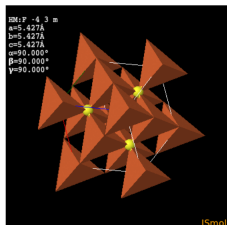COD Wiki
Obtaining COD
Querying COD
Citing COD
COD Mirrors
Advices to donators
Useful links

### Information card for entry 1525302

1525301 << **1525302** >> 1525303

**Preview**



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

JSmol

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 persistence

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

# Raman spectroscopy data

## The Raman Open Database

http://solsa.crystallography.net/rod

**Raman Open Database**

**Information card for entry 3500024**

3500023 << 3500024 >> 3500025

**Preview**

*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 vocabulary in CIF DDLm dictionaries

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 $\rightarrow$ 2.x);

# SOLSA project, COD and ROD



COD will be used in SOLSA for:

- mineral identification;
- subsequent data dissemination.

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

# The fun of REST

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).

# 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'
    ```

# 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
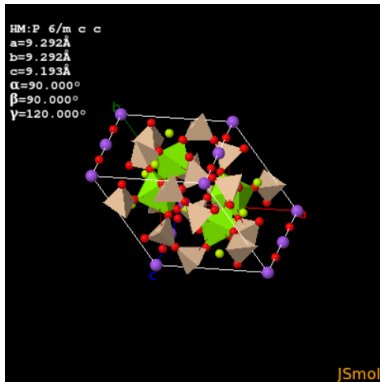Peter Murray-Rust
Miguel Quirós

# Thank you!





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

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

*A path to freedom: GNU → Linux → Ubuntu → MySQL → R → 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.

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

# References II

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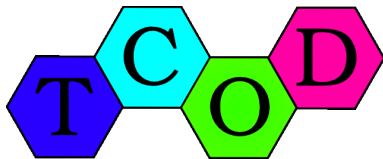Tahoe-lafs distributed storage service in community
network clouds.
*2014 IEEE Fourth International Conference on Big
Data and Cloud Computing.*

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

# Open Crystallographic Databases

COD, TCOD, PCOD, MPOD, ...



http://www.crystallography.net/cod
$> 385\,000$ 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 $> 10^8$?)

*A path to freedom: GNU $\rightarrow$ Linux $\rightarrow$ Ubuntu $\rightarrow$ MySQL $\rightarrow$ R $\rightarrow$ LaTeX $\rightarrow$ TikZ $\rightarrow$ Beamer*

# 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

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

# Hyperspectral image database (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** >> 4060000

**Preview**

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

# 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 client

Tahoe-LAFS
REST
web-API

over HTTP(S)
or (S)FTP

• 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

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

# Tahoe LAFS Grid for SOLSA

## Grid Status

✓ **2 introducers connected**
⊘ **Helper**
  None

**Services**
- Not running storage server
- Not running helper

### Connected to 6 of 6 known storage servers

| Nickname | Connection | Last RX | Version | Available |
|----------|-----------|---------|---------|-----------|
| ✓ **balandis** <br> v0-45zgep2rbv3iwqdn3wdbentg4ioa6iy7zdcuobsxofgp0z | Connected to tcp:172.17.170.119:53026 via tcp | 15h 33m 28s | 1m 5s | tahoe-lafs/1.12.1 | 1867.64GB |
| ✓ **delfinas3** <br> v0-5w2m33nn3h2z4ia7rar5acxpykj3orwmgs42iyoikewiqmrf744za | Connected to tcp:172.17.170.129:51898 via tcp | 15h 33m 28s | 1m 4s | tahoe-lafs/1.12.1 | 469.92GB |
| ✓ **orka** <br> v0-nkd5pakq95ezpwvbva24bxlulyhcwb4itjj3vrrmrt0wq72lrvw4a | Connected to tcp:172.17.170.122:47977 via tcp | 15h 33m 28s | 1m 4s | tahoe-lafs/1.12.1 | 2965.21GB |
| ✓ **stumbras** <br> v0-rezd9j7g9i5unw3mn54ts33njxp3xmgigdigp7czylmx24ik4q | Connected to tcp:172.17.170.121:47082 via tcp | 15h 33m 28s | 1m 4s | tahoe-lafs/1.12.1 | 2965.21GB |
| ✓ **delfinas** <br> v0-ik7ylrn2pnrpt2o2vul6giux2arlsobm2zbxzhwtbxnfnm3xrq | Connected to tcp:172.17.170.129:52200 via tcp | 15h 33m 28s | 1m 4s | tahoe-lafs/1.12.1 | 466.02GB |
| ✓ **delfinas2** <br> v0-smpbw4pzctzu7demrkuuirpd3syypdoaktdo5ppga7ylk7seerrnuja | Connected to tcp:172.17.170.129:34498 via tcp | 15h 33m 28s | 1m 4s | tahoe-lafs/1.12.1 | 469.92GB |

### Connected to 2 of 2 introducers

| | Connection | | Last RX |
|--|-----------|--|---------|
| ✓ | Connected to tcp:172.17.170.121:54295 via tcp | 15h 34m 10s | 1m 29s |
| ✓ | Connected to tcp:172.17.170.122:57127 via tcp | 15h 34m 12s | 1m 47s |

*Tahoe-LAFS for SOLSA set up by Erikas Raginis*

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

# HOD files on the Tahoe LAFS grid

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

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 REST API

- Agreed upon in the 2016 Leiden CECAM workshop;
- Suitable for all structural and QM databases.

https://github.com/Materials-Consortia/API

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

```
(* The top-level 'filter' rule: *)
Filter = Keyword, Expression;
(* Keywords *)
Keyword = "filter=" ;
(* Values *)
Value = Identifier | Number | String ;
(* ... some token definitions skipped for brevity ... *)
(* Expressions *)
Expression = Term, [Spaces], [ OR, [Spaces], Expression ] ;
Term = Comparison, [Spaces], [ AND, [Spaces], Term ] ;
(* Operator Comparison operator tokens: *)
Operator = '<', [ '=' ] | '>', [ '=' ] | '=' | '!', '=' ;
Comparison = Value, [Spaces], Operator, [Spaces], Value |
             NOT, [Spaces], Comparison |
             '(', [Spaces], Expression, [Spaces], ')' ;
```

# Schemas for return data

Schemas allow to:

- formally agree on what is right and wrong;
- validate program outputs and documents automatically.

```
"query": {
    "type": "object",
    "properties": {
        "representation": { "type": "string" },
        "api_version": { "type": "string" },
        "time_stamp": { "type": "string" },
        "data_returned": { "type": "integer" },
        "data_available": { "type": "integer" },
        "last_id": { "type": "string" }
    },
    "required": [ "representation", "api_version",
                  "time_stamp" ]
},
```

# API query examples

http://crystallography.net/cod/optimade/structures?filter=elements="Si,O"ANDnelements=2&limit=1

```
{
  "resource": {
    "base_url": "http://www.crystallography.net/cod/optimade/v1.0-alpha.1/"
  },
  "query": {
    "api_version": "v1.0-alpha.1",
    "data_returned": 1,
    "representation": "/structures?filter=elements=\"Si,O\"ANDnelements=2&limit=1",
    "last_id": "1010921",
    "time_stamp": "2017-04-06T05:46:50Z",
    "implementation": {
      "maintainer": {
        "email": "cod-bugs@ibt.lt"
      },
      "title": "Crystallography Open Database",
      "version": "v1.0-alpha.11",
      "source_url": "svn://crystallography.net/cod/trunk/cod/cgi-bin/optimade.pl@194653"
    },
    "data_available": 344
  },
  "data": [
    {
      "last_modified": "2017-02-28T05:33:56Z",
      "properties": {
        "formula": "O2 Si"
      },
      "url": "http://www.crystallography.net/cod/1010921.cif",
      "immutable_id": "http://www.crystallography.net/cod/1010921.cif@130149",
```

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

# Common pattern of self-describing data definitions

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