



Australian Government



# What is a dataset?

---

James Hester

Australian Centre for Neutron Scattering

Science. Ingenuity. Sustainability.



The Future of Research Communications and e-Scholarship

- ▶ (F)indable: Keywords, semantic web, registries, DOIs
- ▶ (A)ccessible: DOIs, protocols
- ▶ (I)nteroperable: Standards for data description
- ▶ (R)eusable: Standards

# The "data library"

You've found the book: now what?



# What is a dataset?

- ▶ All the data required to support a (published) result?
- ▶ All the data collected during an experimental run?
- ▶ All the data necessary to assess a theoretical calculation?
- ▶ ...

*A dataset is the complete data required for a particular purpose*

# How do we describe datasets to software?

- ▶ Do not necessarily coincide with the contents of a DOI
  - ▶ A data collection split over several DOIs
  - ▶ Several datasets in a DOI data dump
- ▶ File formats vary and can be mixed together

# A motivating example: CheckCIF for raw data



Check raw data files for

- ▶ consistent metadata
- ▶ sufficient metadata
- ▶ agreement with results

But ...

Raw data diffraction files come in over 200 different styles (V. Minor)

# Solution

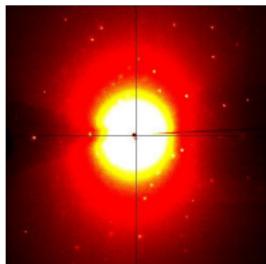
A collection of data files = a collection of tables

(see Hester, J.R. "A robust, format-agnostic scientific data transfer framework", (2016) *Data Science Journal* 15 12:1-17)

*A dataset specification = a set of column headings*

# First: map file contents to tabular form

```
{  
HEADER_BYTES= 512;  
DIM=2;  
BYTE_ORDER=little_endian;  
TYPE=unsigned_short;  
SIZE1=2048;  
SIZE2=2048;  
PIXEL_SIZE=0.1024;  
BIN=2x2;  
ADC=fast;  
DETECTOR_SN=457;  
TIME=2.000000;  
DISTANCE=68.042500;  
PHI=4.000000;  
OSC_START=4.000000;  
OSC_RANGE=2.000000;  
WAVELENGTH=0.710894;  
BEAM_CENTER_X=105.134680;  
BEAM_CENTER_Y=104.653206;  
ACC_TIME=3156;  
CREV=1;  
BIN_TYPE=HW;
```



## ADSC (SMV) format

- ▶ One file per oscillation step
- ▶ Header + data

# ADSC files as a collection of tables

| Frame_id | Time | Detector_SN | Wavelength |
|----------|------|-------------|------------|
| kt_001   | 2.0  | 457         | 0.710894   |
| kt_002   | 2.0  | 457         | 0.710894   |
| kt_003   | 2.0  | 457         | 0.710894   |

| Frame_id | Axis | Displacement | Angle |
|----------|------|--------------|-------|
| kt_001   | phi  | .            | 0.0   |
| kt_001   | dist | 68.042       | .     |
| kt_002   | phi  | .            | 2.0   |
| kt_002   | dist | 68.042       | .     |
| kt_003   | phi  | .            | 4.0   |
| kt_003   | dist | 68.042       | .     |

# Grouping data files

| Frame_id | Time | Detector_SN | Wavelength | Scan_id |
|----------|------|-------------|------------|---------|
| kt_001   | 2.0  | 457         | 0.710894   | 1       |
| kt_002   | 2.0  | 457         | 0.710894   | 1       |
| kt_003   | 2.0  | 457         | 0.710894   | 1       |

| Frame_id | Axis | Displacement | Angle |
|----------|------|--------------|-------|
| kt_001   | phi  | .            | 0.0   |
| kt_001   | dist | 68.042       | .     |
| kt_002   | phi  | .            | 2.0   |
| kt_002   | dist | 68.042       | .     |
| kt_003   | phi  | .            | 4.0   |
| kt_003   | dist | 68.042       | .     |

# imgCIF: A "master" image ontology

International Tables for Crystallography (2006). Vol. G, Chapter 4.6, pp. 444–450.

## 4.6. Image dictionary (imgCIF)

BY A. P. HAMMERSLEY, H. J. BERNSTEIN AND J. D. WESTERROCK

This is version 1.3.2 of the Image CIF dictionary (imgCIF) and crystallographic binary file dictionary (CBF) extending the structures contained in the International Tables for Crystallography, as described in Chapter 3.7. See also Chapter 2.3 for a description of the CBF format and Chapter 5.6 for discussion of a software library for reading and writing CBF files.

There are three category groups in this dictionary:  
array\_group contains categories that describe array data;  
array\_group contains categories that describe the axes; and  
diffraction\_group contains categories that describe details of the diffraction experiment.

```
imgcif_1.3.2
    <CIF-IMAGE-DICTIONARY>
        <CONSTANT-DEFINITION>
            <CONSTANT-APPLICATION-STRUCTURE>
                <CONSTANT-SCALAR-STRUCTURE>
                    <CONSTANT-FLOATING-POINT-STRUCTURE>
                        <CONSTANT-BYTE-APPLICATION-STRUCTURE>
                            <CONSTANT-BINARY-STRUCTURE>
                                <CONSTANT-BEAN>
                                    <CONSTANT-STRUCTURE>
                                        <CONSTANT-STRUCTURE-NAME>
                                            <CONSTANT-STRUCTURE-ID>
                                                <CONSTANT-MIME>
                                                    <CONSTANT-MIME-ID>
                                                        <CONSTANT-MIME-NAME>
                                                            <CONSTANT-MIME-DESCRIPTION>
                                                                <CONSTANT-MIME-STRUCTURE>
                                                                <CONSTANT-MIME-STRUCTURE-ID>
```

```
<CONSTANT-MIME-STRUCTURE>
<CONSTANT-MIME-STRUCTURE-ID>>
    <CONSTANT-MIME-NAME>
    <CONSTANT-MIME-DESCRIPTION>
        <CONSTANT-ARRAY-STRUCTURE>
            <CONSTANT-ARRAY-NAME>
                <CONSTANT-ARRAY-ID>
                    <CONSTANT-ARRAY-ITEMS>
                        <CONSTANT-ARRAY-ITEM>
                            <CONSTANT-ARRAY-ITEM-ID>
                                <CONSTANT-ARRAY-ITEM-NAME>
                                    <CONSTANT-ARRAY-ITEM-ID>
                                    <CONSTANT-ARRAY-ITEM-NAME>
                                    <CONSTANT-ARRAY-ITEM-DESCRIPTION>
                                    <CONSTANT-ARRAY-ITEM-ITEMS>
                                        <CONSTANT-ARRAY-ITEM-ITEM>
                                            <CONSTANT-ARRAY-ITEM-ITEM-ID>
                                                <CONSTANT-ARRAY-ITEM-ITEM-NAME>
                                                <CONSTANT-ARRAY-ITEM-ITEM-ID>
                                                <CONSTANT-ARRAY-ITEM-ITEM-NAME>
                                                <CONSTANT-ARRAY-ITEM-ITEM-DESCRIPTION>
```

The example shows two binary data blocks. The first one was composed by the CBF encoder and the second one was composed by the imgCIF encoder. The class name 'X' is on the data items to be encoded. In both cases X is identified as 'array'. The first data item has the array identifier 'array\_data\_array\_id'. The second data item has the array identifier 'array\_data\_binary\_id'. In the case of 'X', which then requires eight hexadecimally digit bytes. The third data item has the array identifier 'array\_data\_structured\_id'. The fourth data item has the array identifier 'array\_data\_structured\_id'. Additionally, the reference '...' could have been used for the reference 'array'. Both encoders support the structured representation and the structured representation is enabled for binary data by specifying the codecs and CBFs encoding. Note that the structured representation is not supported by the CBF encoder.

```
<CONSTANT-ARRAY-ITEM-ITEM-ITEMS>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-ID>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-NAME>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-ID>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-NAME>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-DESCRIPTION>
<CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEMS>
    <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM>
        <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ID>
            <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-NAME>
            <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ID>
            <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-NAME>
            <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-DESCRIPTION>
            <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEMS>
                <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM>
                    <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM-ID>
                        <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM-NAME>
                        <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM-ID>
                        <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM-NAME>
                        <CONSTANT-ARRAY-ITEM-ITEM-ITEM-ITEM-ITEM-DESCRIPTION>
```

The example shows two data blocks. The first one was composed by the CBF encoder and the second one was composed by the imgCIF encoder. The constant application structure is a variant of the Multipurpose Internet Mail Extensions (MIME) specified in RFC 2045-2049 by N. Freed et al., the header name is 'Content-Type'. The value of 'Content-Type' in a CBF is '<CIF-IMAGE-DICTIONARY>' (including the required initial '<'). The value of 'Content-Type' in a CBF is 'application/cif-dict' and is RFC 2045: 'application/cif-dict' is recommended. If an exact structure was composed, the composition should be specified by the name of the structure 'Content-Type' or 'Content-Transfer-Encoding' or 'Content-MIME-Type'. The Content-MIME-Header 'Content-MIME-Type' is 'application/cif-dict' and is RFC 2045: 'application/cif-dict' is recommended. The Content-Transfer-Encoding may be 'BASE64', 'Quoted-Printable', 'X-BASE64', 'X-QUOTED-PRINTABLE', 'X-MIME-ASN1' or 'X-MIME-AES'. For a CBF, the octal, decimal and hexadecimal transfer encodings are not recommended for binary data. The CBF encoder adds an empty line at the end of each block of binary data before the empty line terminating the header. In a CBF, the raw binary data begin after the empty line terminating the header. The octal, decimal and hexadecimal transfer encodings are not recommended for binary data. The CBF encoder adds an empty line at the end of each block of binary data before the empty line terminating the header. In a CBF, the raw binary data begin after the empty line terminating the header.

Octal Hexadecimal Decimal Purpose  
0 00 20 0x00 0x00 0x00 page break  
1 A 2B 0x1A 0x2B Cif-Z: strip listing, MS-DOS  
2 04 08 0x04 0x08 Cif-D: strip listing, UNIX  
3 DF 23 0x0D 0x15 Cif-D: strip listing

None of these octets are included in the calculation of the message size or in the calculation of the message digest. The X-Binary-Size

Abbreviations. ANTHONY P. HAMMERSLEY, CNRS-ESPCI, Paris, France. HILBERT J. BERNSTEIN, Department of Mathematics and Computer Science, Bar-Ilan University, Ramat-Gan, Israel. JOSEPH D. WESTERROCK, IBM T.J. Watson Research Center, Yorktown Heights, NY 10598, USA. JOHN D. WETTERSTROM, Dept. of Materials Science and Engineering and Chemical Biology, The State University of New Jersey, Rutgers, New Brunswick, NJ 08854-8087, USA.

Copyright © 2006 International Union of Crystallography

444 International Tables for Crystallography (2006). Vol. G, Chapter 4.6, pp. 444–450.

# Ontologies

A collection of definitions for data names



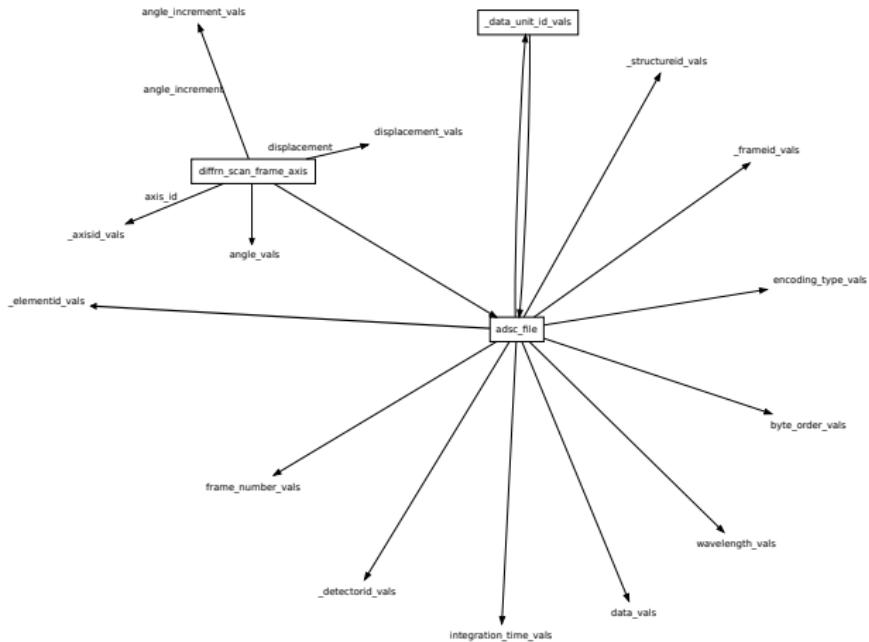
# Category theory

"Objects" connected by "morphisms"

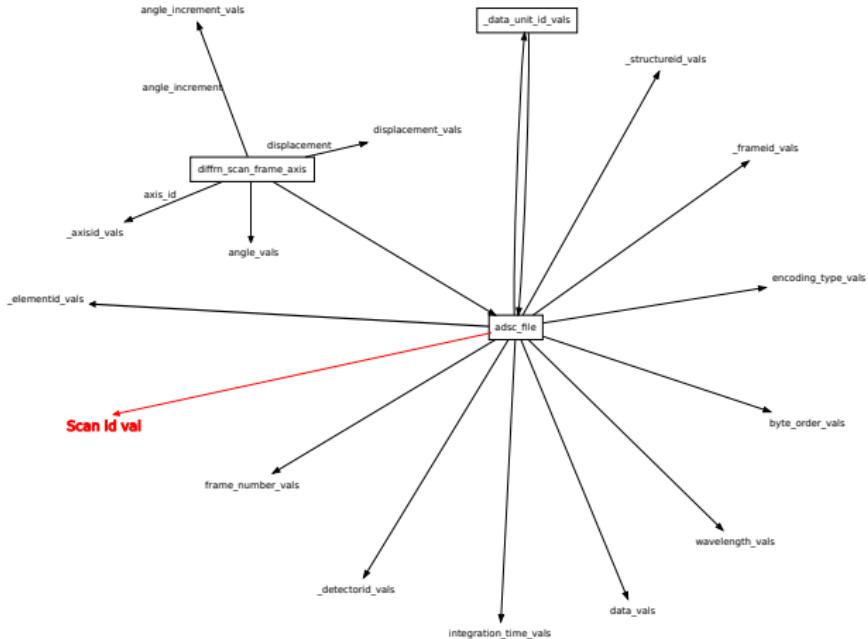


A collection of tables (a relational database) = a category over sets and functions

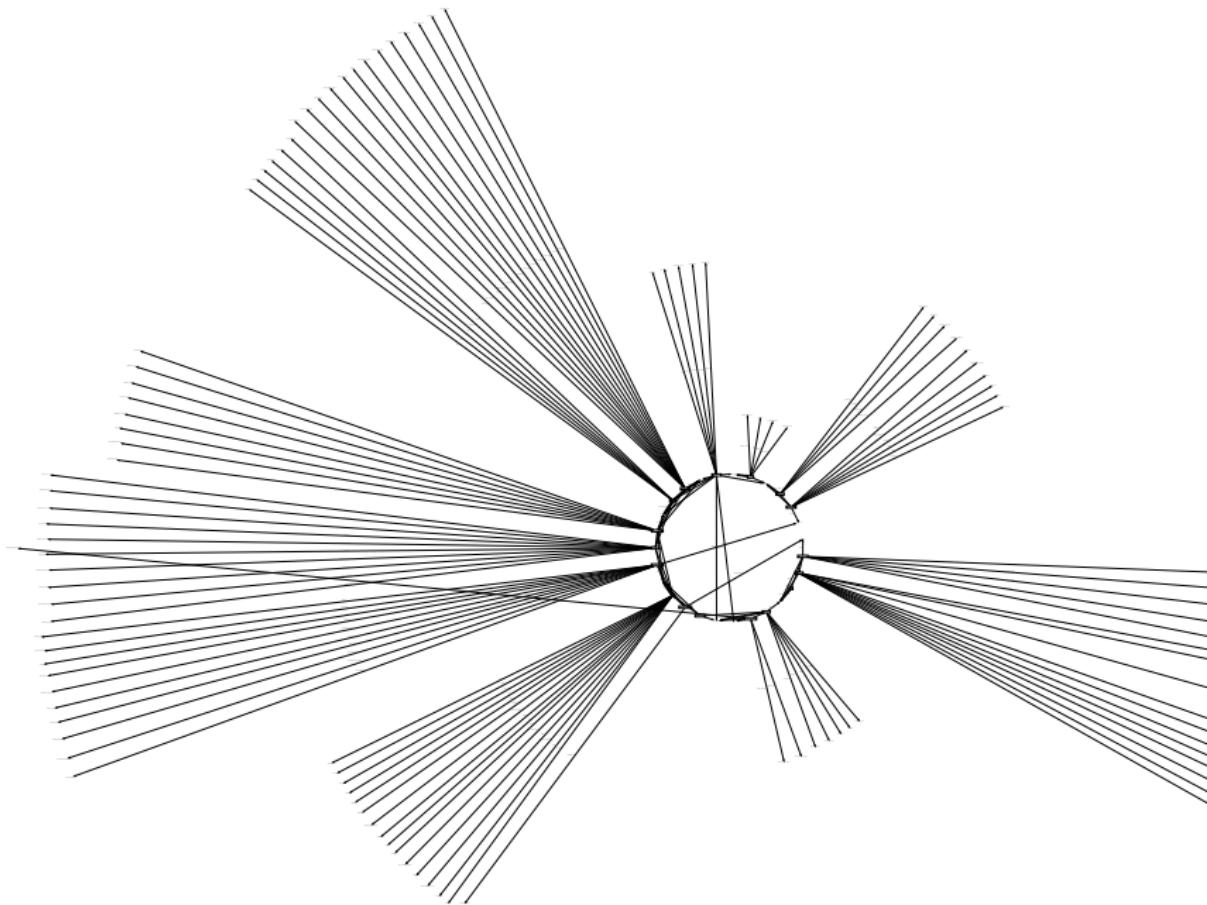
# ADSC file collection as a category



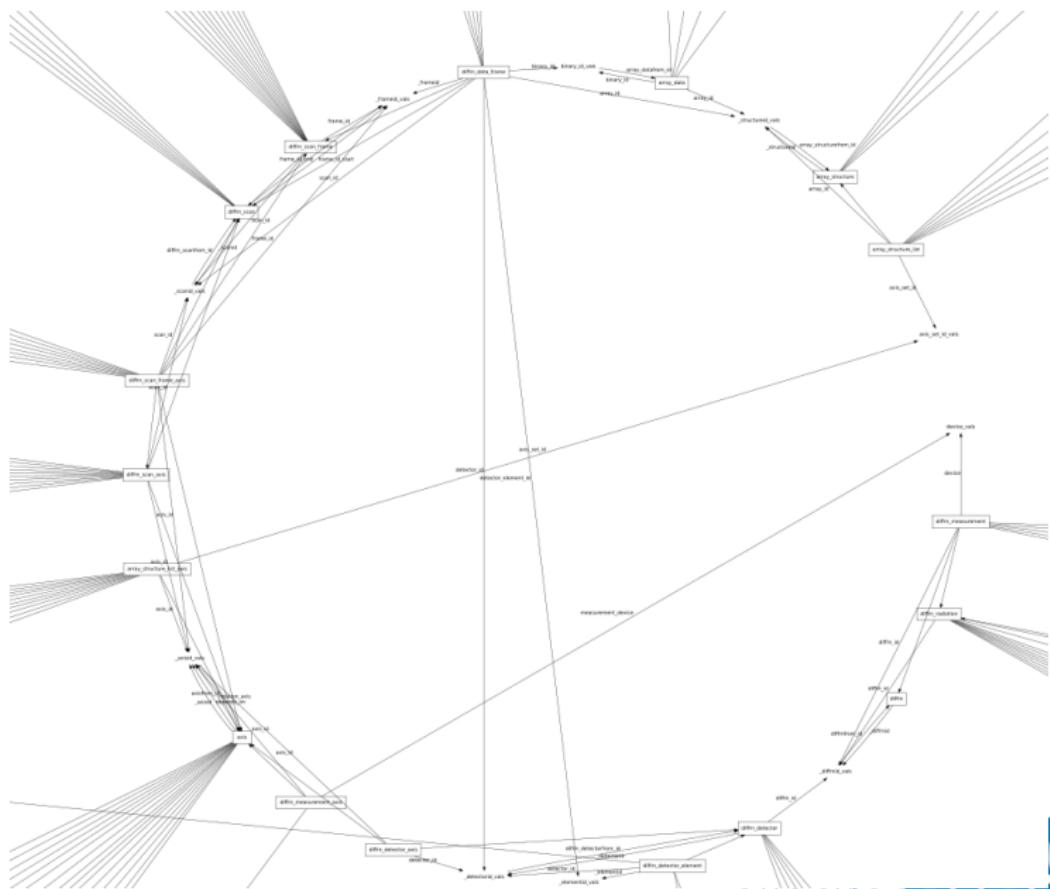
# Grouping into scans



# imgCIF diagrammed



# imgCIF, zoomed in



# Mappings between ontologies (Functors)

Objects go to objects, paths go to paths

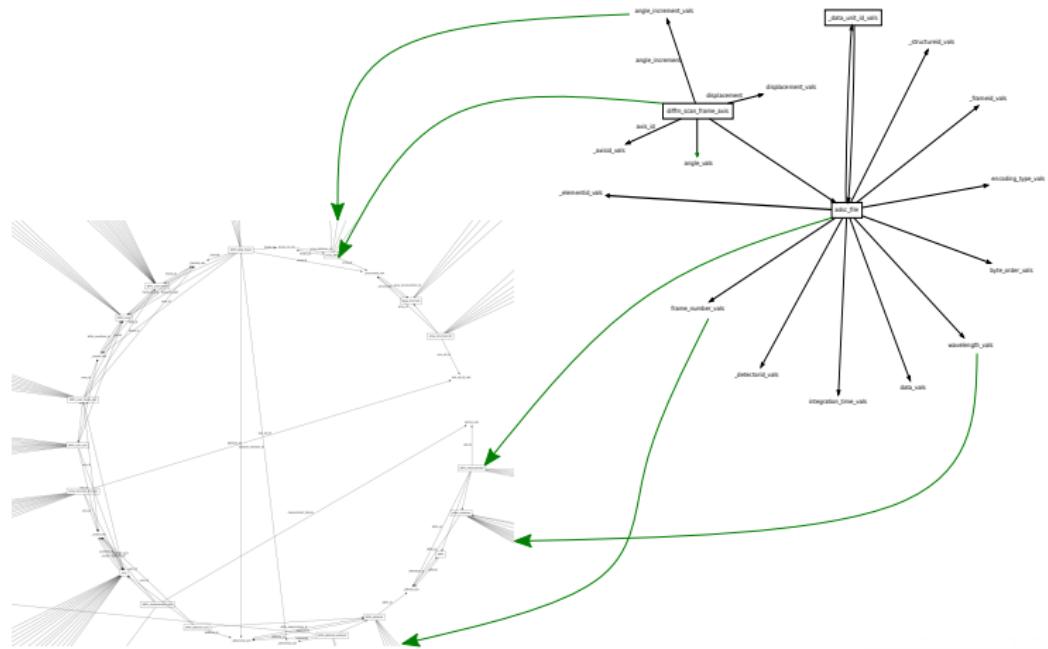
Given a functor, automatic data migration **between categories** possible:

- ▶ Left and right pushforward functors
  - ▶ Equivalent to database joins and unions
- ▶ Machine-actionable
- ▶ Fills in the gaps

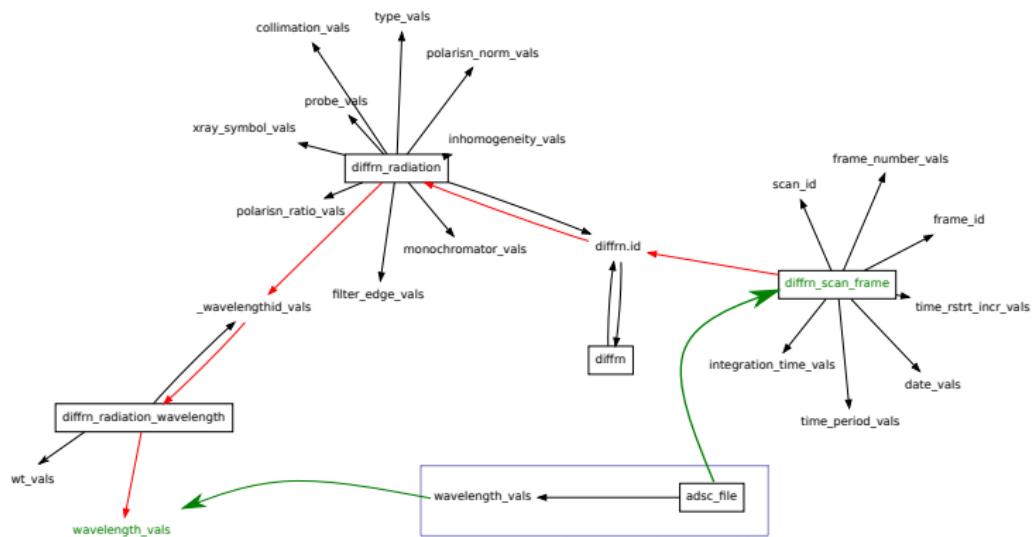
Spivak, D. I. "Functionial Data Migration" *Information and Computation* 217 (2012) 31-51

Useful also unpublished earlier version <http://math.mit.edu/~dspivak/informatics/FunctionialDataMigration.pdf>

# Functor from ADSC category to imgCIF category



# In detail



# Summary

- ▶ Tabular (categorical) representation is basis for interoperability standards
  - ▶ Machine-actionable descriptions of arbitrary formats
  - ▶ Links to community ontologies
- ▶ Need to specify only equivalences between tables and columns

# A "virtual depository"

1. Provide one or more DOIs
2. System determines and lists file types found at those DOIs
3. Assign format specifiers to each file type

Later....

1. Request a dataname
2. System returns values for this dataname
3. More sophisticated: return dataname corresponding to a given key

# Outlook

- ▶ Proof of concept almost there
- ▶ Scaling up requires fast software
- ▶ Julia allows leveraging of existing C, Fortran libraries