

FAIR Data – why and how

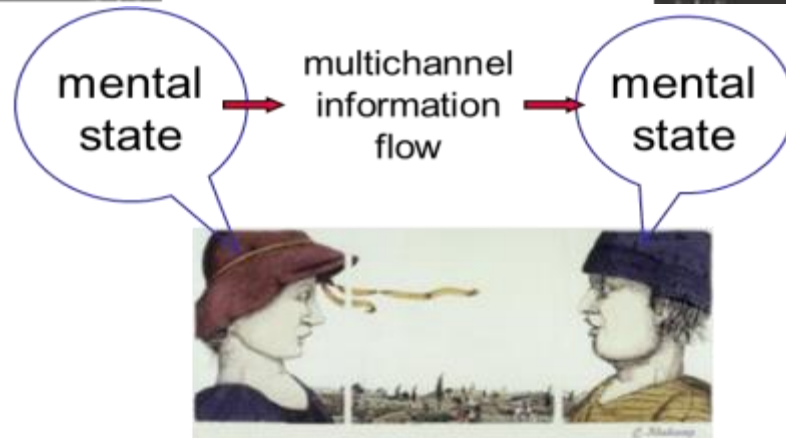
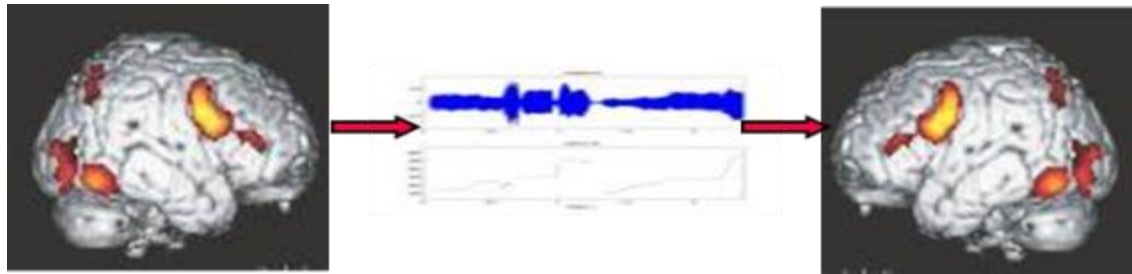
Peter Wittenburg

FAIR DIGITAL OBJECTS  FORUM

1. Short bio
2. About metadata, persistent IDs and the FAIR Principles
3. State of data
4. Globally integrated dataspace and FAIR Digital Objects

Short Bio

Responsible for technology and methodology in a Max Planck Institute for understanding the human brain



- Modern science is data driven
- Applied advanced pattern recognition using data hungry algorithms
- Started building data infrastructures to improve access to data from other institutes
- Responsible for DOBES, CLARIN and EUDAT research/data infrastructures
- Started Research Data Alliance 2013
- Started FAIR DO Forum 2019

Library Scene

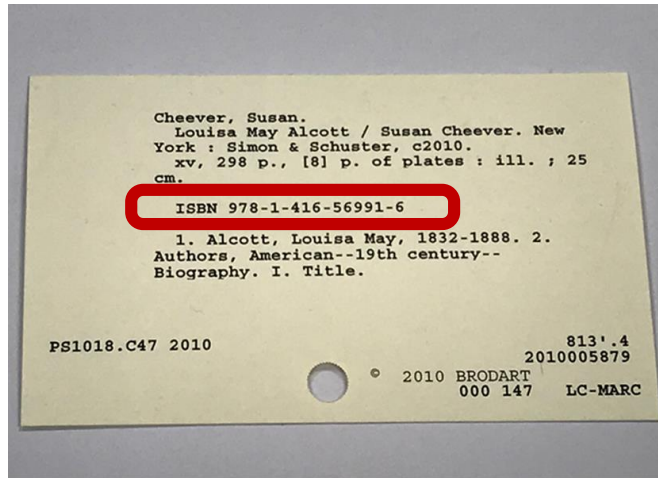
I am afraid I can't help you. Do you have the ISBN number, author name, etc.?



I need the book that was in row 20/shelf 11/layer3 about 5 years ago.

Of course, the place in a shelf is NOT a persistent identifier of a book. The library might have undergone a change, might have several copies, etc.

Library PIDs and Metadata



Transformation to
an all digital world



EUDAT B2Share Repository

SEARCH

HELP COMMUNITIES UPLOAD CONTACT

Login

RECORDS - 4E8AC36CDD343DA81FD9E83E72805A0

Latest version - Mar 2, 2018

Common Patterns in Revolutionary Infrastructures and Data

by Peter Wittenburg, George Strawn;

Mar 2, 2018

	Views	File Downloads
	46	28

Description: Large infrastructures follow certain patterns during their evolution. The paper looks into three examples (electrification, Internet, WorldWide Web), extracts some reoccurring patterns and finally compares it with th state and potential of the data domain. Summarising it states that we seem to be close to a convergence step, but that we lack the agreement yet on a fairly simple concept all could agree with and that can serve as commodity to build upon.

Disciplines: 4.1.5 → Computer sciences → Data structures

Keywords: infrastructure evolution, data infrastructure, digital

DOI: [10.23728/b2share.4e8ac36cdd343da81fd9e83e72805a0](https://doi.org/10.23728/b2share.4e8ac36cdd343da81fd9e83e72805a0)

PID: [11304/34c5c223-adf9-49fa-9afe-968a7a932970](https://nbn-resolving.org/urn:nbn:de:hbz:5:1-11304-34c5c223-adf9-49fa-9afe-968a7a932970)

A book has a Persistent ID (ISBN)
and a Metadata Card.

It will direct you where you can find
the location of the printed copy.

Metadata according to a simple
schema (structure & vocabulary)

A digital object has a DOI/Handle
as Persistent ID.

A URL is a location in a store!

Many Types of Metadata

Digital Calibration Certificate

1.: Administrative Data

- regulated

2.: Results of the Calibration

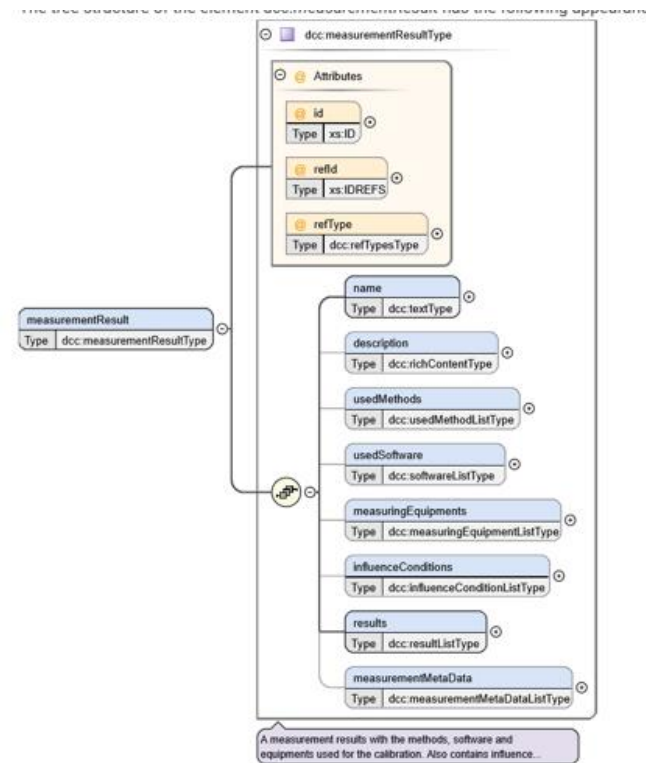
- regulated:
 - $Y = y \pm U(k) [S]$
- not regulated:
 - individual data
 - ...

3.: Comments

- not regulated

4.: Document

- human readable



Extract from the XML schema: type definition "measurementResultListType"

- DCC is a typical complex metadata schema
- Describes the act of a calibration to document quality etc.

Many Types of Metadata

Digital Calibration Certificate

1.: Administrative Data

- regulated

2.: Results of the Calibration

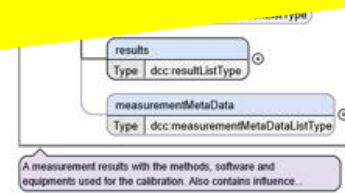
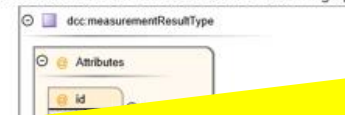
- regulated:
 - $Y = y \pm U(k)$ [SI]
- not regulated:
 - individual data
 - ...

3.: Comments

- not regulated

Metadata and Persistent Identifiers are well-established Concepts.
Librarians have shown their importance for many activities.
The digital domain is more diverse and complex.
Principles seem to be eternal.

- DCC is a complex document
- DCC is a part of a document
- DCC is a specification
- DCC requires access rights and requires signatures



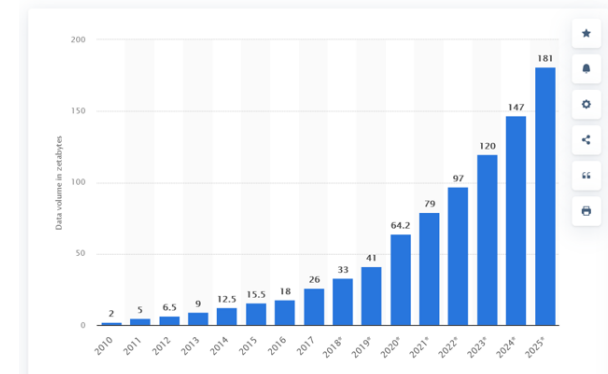
Extract from the XML schema: type definition "measurementResultListType"

Books and Data

- Libraries are about books, papers and media (publications)
 - there are about 150 million books in total
 - books are human readable and in libraries (findable, accessible, interoperable, reusable)
- Research repositories are about reusable data in varying circumstances
 - Our small MPI data repository had already ~1 million digital data objects in 2012
 - Metadata: Searching for data with same brain diseases, same age, etc.
 - Climate simulation people create about 100.000 DOs in a simulation run
- The number of Digital Objects cannot be counted anymore
 - Quantity requires new mechanisms
- Of course, want to access all research artefacts
 - Open Science in 16xx: make results available to peers
 - Open Science in 202x: make research artefacts available

Technology & Telecommunications · Telecommunications

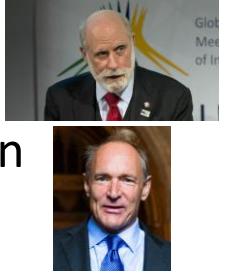
Volume of data/information created, captured, copied, and 2010 to 2020, with forecasts from 2021 to 2025
(in zettabytes)



State of Data



- Surveys: 80% of time and money in data-projects is wasted with data munging
- P. Heidorn (U California): in 20 years > 80% of our data will be lost
- V. CERF (Internet pioneer): we are risking the Digital Dark Age
- R. Kahn (Internet pioneer): we lack a persistent, stable & global data management solution
- T. B. Lee (Web pioneer): we need a radical change – a new architecture



- Digital data to become findable
 - metadata is required - supporting domain vocabularies
- Digital data to become accessible
 - Needs a general access protocol
- Digital data to become interoperable (not only human readable)
 - software is required – needs structure and semantic definitions
- Digital data to become reusable
 - Metadata is mostly poor and rights are not clear



FAIR \neq Open Data

State of Data

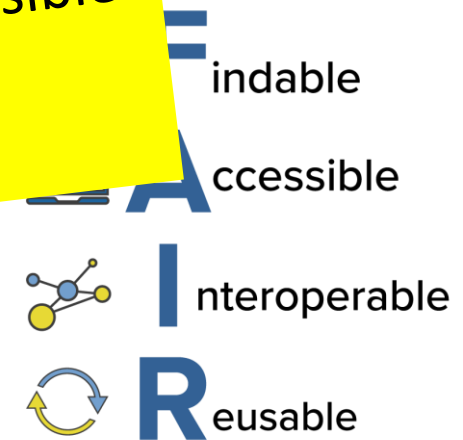


- Surveys: 80% of time and money in data-projects is wasted with data munging
- P. Heidorn (U California): in 20 years > 80% of our data will be lost
- V. CERF (Internet pioneer): we are risking the Digital Revolution
- R. Kahn (Internet pioneer): we are risking the Digital Revolution
- T. Berners-Lee: we are risking the Digital Revolution



The state of data is bad.
Measures of improvements are urgently required to improve the state of data.
Need to urgently reduce costs to make data projects feasible for many.

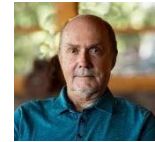
- Digital data to become findable
 - metadata
- Digital data to become accessible
 - Needs
- Digital data to become interoperable (not only human readable)
 - software is required – needs structure and semantic definitions
- Digital data to become reusable
 - Metadata is mostly poor and rights are not clear



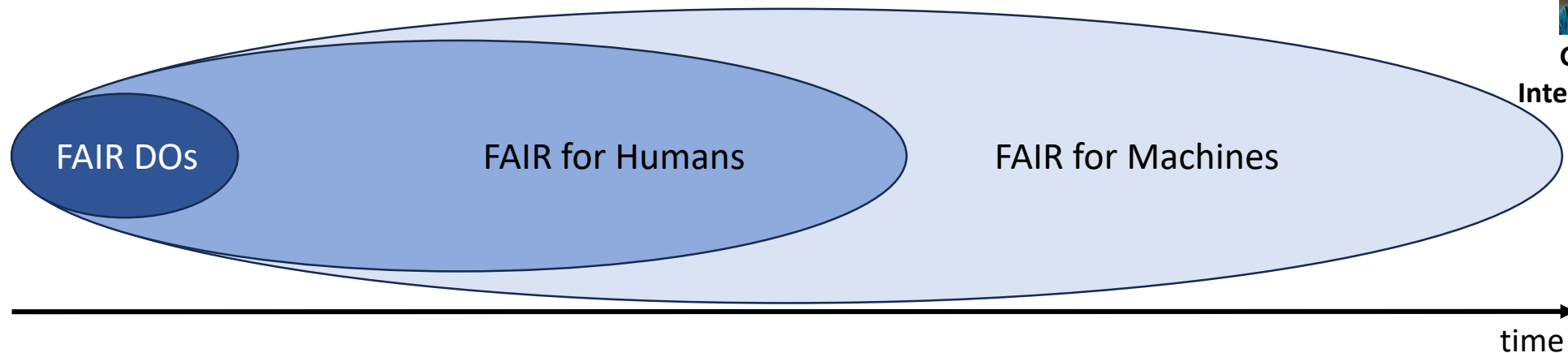
FAIR ≠ Open Data

Feasibility

- 100 % FAIR implies a huge transformation of repository strategies
 - Globally we have thousands of repositories – mostly underfunded
 - This transformation will cost much money and time
 - A fast integration as a first step is, however, possible for far less money

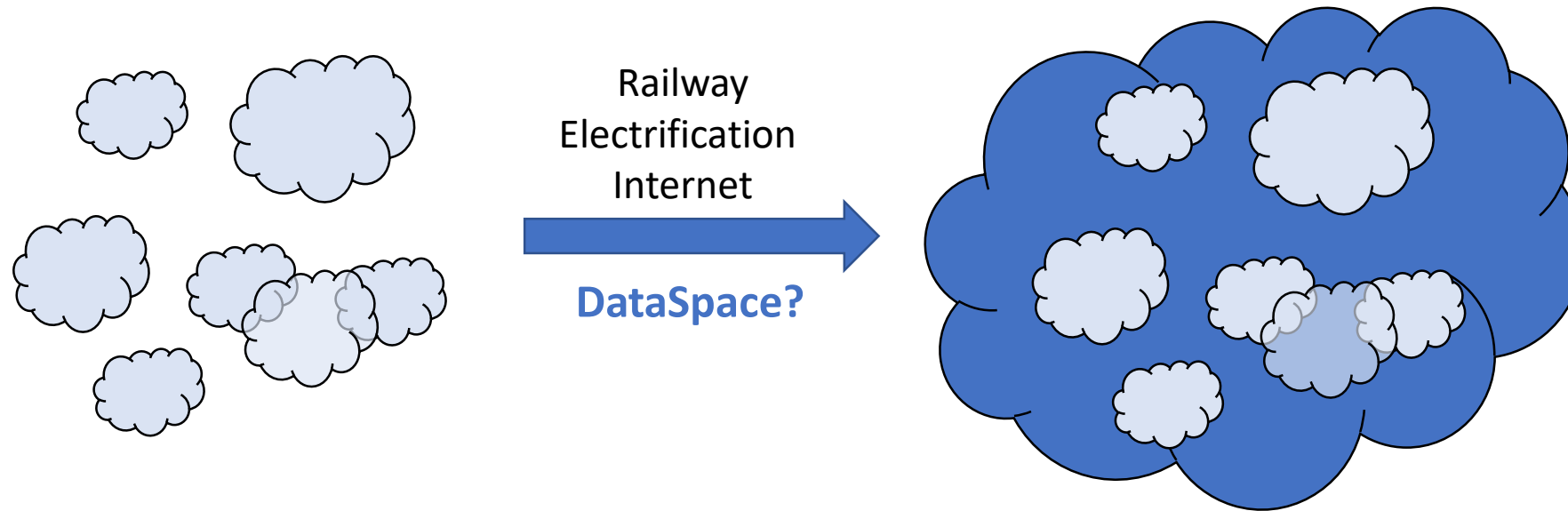


G. Strawn
Internet Pioneer



- For a basic integration we need to repeat the Internet success story
 - Creating a minimal Interop standard to connect all existing and emerging islands
 - Not trying to change everything at once

Globally Integrated Dataspace



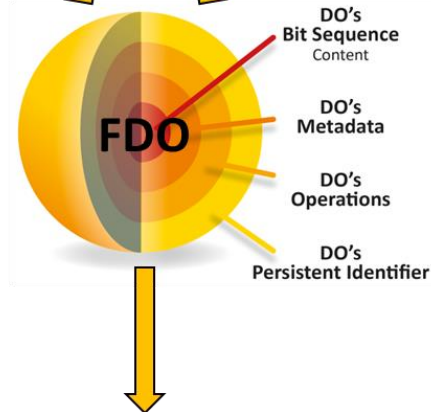
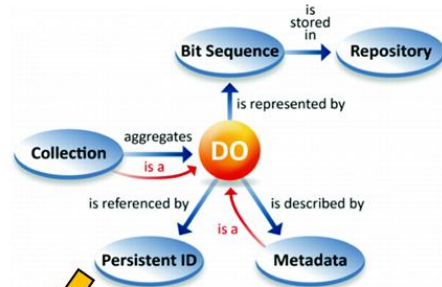
- Needs to be based on FAIRness and increase TRUST
 - Mechanisms to be based on Responsibility, Accountability and Persistency
 - Mechanisms should be distributed and support global governance

FAIR Digital Objects similar to Internet Basics

FAIR Principles

Findable 
Accessible 
Interoperable 
Reusable 

Digital Objects



FDOs

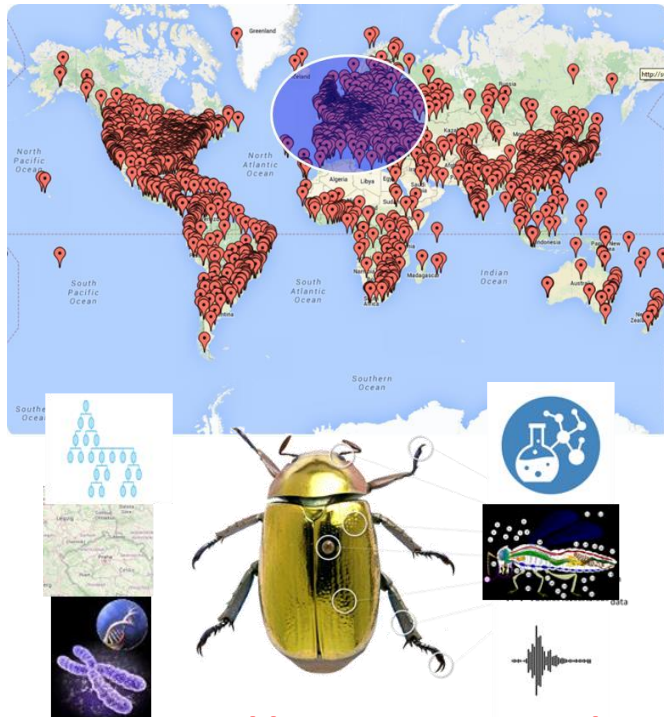
- are persistent atomic entities of FAIR information identified by a PID
- follow minimal specifications about persistently bundling all information to enable processing of an object
- do not require to have FAIR repositories right away
- increase the protection of data where necessary

FDO Forum

- Independent
- No patents etc. involved
- participate

2023: FDO Requirements Specifications ready
March 2024: FDO Summit with many applications

FDOs to structure the complex Data Landscape



<https://www.dissco.eu/>

20.5000.1025/xyxy123 ODType1803

PID DOIType

name: Conus Hughmorrisoni
size: 20.4 cm
depth: 9-11 m
locality: Kavieng Lagoon
country: Papua New Guinea
collector: etc. etc.

Metadata *

Op 123/123 (return MIDIS info)
Op 123/124 (update MIDIS info)
Op 123/125 (present hi-res images)
Op 123/556 (return gathering)
Op 123/098 (return DNA)

Invoke image serving
Link to gathering site and other event info

Local id: MNHN-IM-2013-53462

Related objects
Hi-res image 1
Hi-res image 2
Hi-res image N

Link to DNA Data

- global Biodiversity-Infrastructure
 - 120 Natural Science Institutes in EU
 - 1000 collections, 3 billion objects, 2 million standards
 - trillions of relations (storing knowledge)
 - Scattered information about specimen
- Task1: build on integrated dataspace and virtually bundle all information belonging to one specimen
- Task2: enable the execution of specific software on the digital twin of an object (which is an FDO)

FDOs to structure the complex Data Landscape



- global...

FDOs can be used NOW to build a stable global data infrastructure.
FDOs can be used NOW to structure complex dataspaces and thus create knowledge structures for coming generations.
FDOs can be combined with data protection technology (EDC).
Perhaps: FDOs could form the basis for a global quality infrastructure for metrology.



Thanks a lot for your attention.

FAIR DIGITAL OBJECTS  FORUM

<https://fairdo.org/>