

194.044 Data Stewardship (VO 2,0) 2022S
and
194.045 Data Stewardship (UE 2,0) 2022S



Andreas Rauber



Tomasz Miksa



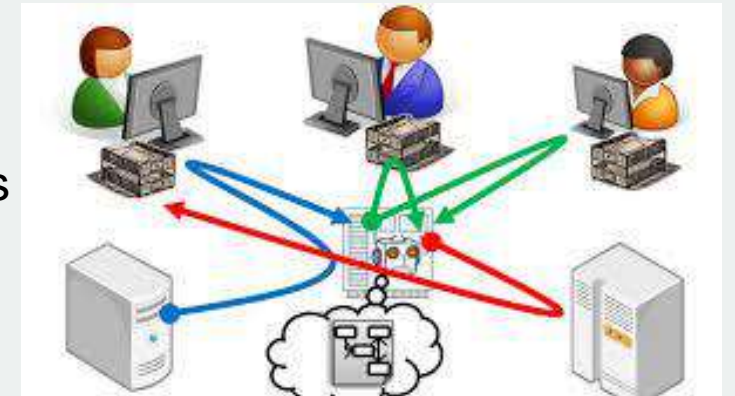
Martin Weise

Agenda

- Organisational
 - What is data stewardship?
 - Overview of lecture topics
 - Dates, exams, materials...
 - Exercise overview
 - Selected projects related to this course
 - (Funded) collaboration possibilities
- Introduction to FAIR principles

Data Science

- Data
 - is fuel for research
 - is the result of processes such as
 - capturing, pre-processing, transformation, integration, analysis
- Research Infrastructures
 - [Human Brain Project](#)
 - neuroscience, computing, and brain-related medicine
 - [Elixir](#)
 - bioinformatics
 - Large Hadron Collider at CERN
 - 300GB per second of raw data from detectors



Researchers trying to reuse data...

Conversation of two researchers

- Can I see your data?
- It's on my USB stick
- Can I have it?
- I have in a box and I have moved recently
- Can I have it?
- I forgot to label the boxes...
- (half a year later)
- Thanks, for the USB. However, I cannot read the hexadecimal file on it. How do I open it?
- You need a special program
- What program?
- ...



Hanson, Karen; Surkis, Alisa; Yacobucci, Karen: Data Sharing and Management Snafu in 3 Short Acts.
<https://doi.org/10.5446/31036>

Data Stewardship (DaSt)

Data Stewardship ~ Data Management ~ Digital Curation

Management of data to ensure its

- FAIR
- Reproducibility
- Auditability
- ...

Must be addressed at different levels and is interdisciplinary

- Technical
- Organisational
- Legal

Data Stewardship is needed in all domains

- E.g. Data Stewards must have domain specific knowledge and data management skills

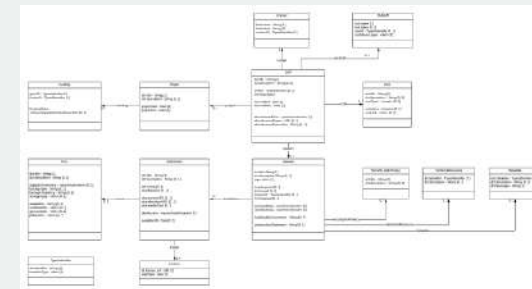
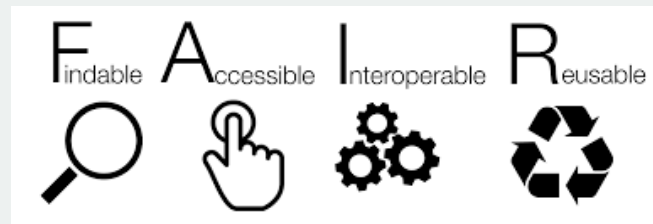


(c) Austrian Airlines / CC-BY-SA

Overview of Lecture Topics

Data Management

- Plans, policies, support services, stewards, etc.
- Legal frameworks, roles and responsibilities, etc.



Overview of Lecture Topics

Data repositories

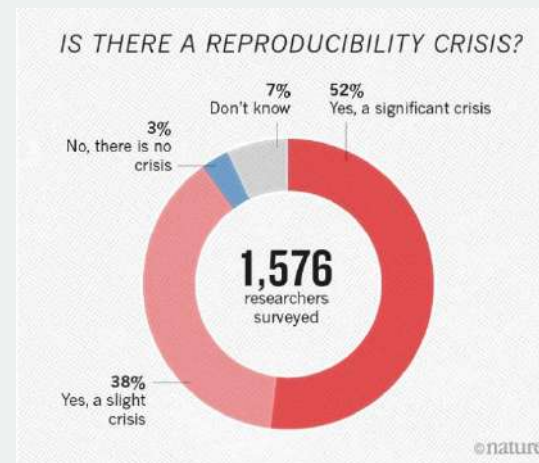
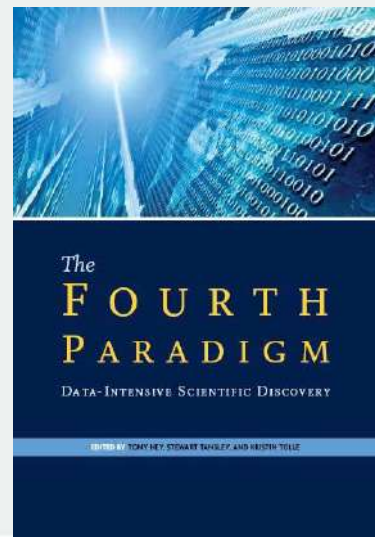
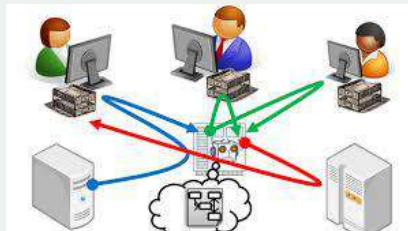
- Software architecture, interfaces, interoperability, data models, standards, ontologies, etc.
- Certification, processes, management, etc.



Overview of Lecture Topics

Reproducibility

- Software execution tracing, dependencies, software packaging, etc.
- Data identification and citation, provenance, etc.



RDA

Scalable dynamic-data
Citation Methodology

The Challenge:
Supporting accurate citation of data subjected to change, for the efficient processing of data and linking from publications.

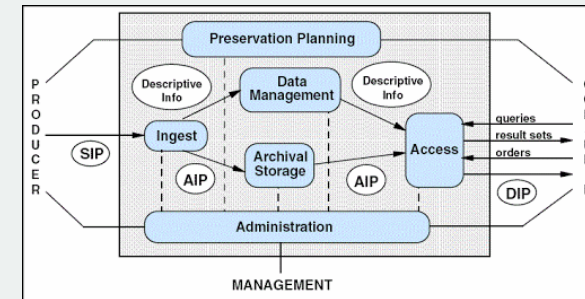
What is the solution?
The Dynamic Data Citation Working Group developed a simple, scalable mechanism that allows the precise, machine-actionable identification of arbitrary sub selections of data at a given point in time irrespective of any subsequent addition, deletion or modification.

Produced by: **Data Citation (WGDC) WG**
<https://rd-alliance.org/groups/data-citation-wg.html>

Overview of Lecture Topics

Digital Preservation

- migration, emulation, planning, risk management, etc.



Lecture Goals

As a student/researcher/employee

- Better design and organise your work with data (and code)
 - Select proper tools and software
 - Better design data transformation processes
 - Make aware decisions on experiment/tool/process design with respect to data
- Make your results trustworthy, auditable, reproducible, FAIR...

As a data steward/consultant/architect

- Have a good understanding of all components needed to establish systems/workflows/practices at an organisation

Lecture Dates

07 March 2022 TM: Intro

14 March 2022 AR: DP

21 March 2022 AR: OAIS

28 March 2022 TM: DMPs

4 April 2022 AR: Data Citation

11 April 2022 #Easter

18 April 2022 #Easter

25 April 2022 TM: FAIR

02 May 2022 TM: maDMPs

09 May 2022 TM: Repositories 1

16 May 2022 TM: Repositories 2

23 May 2022 TM: Developing RDM Services + Certification

30 May 2022 VD: Legal aspects and MW: Secure Data Infrastructures

13 June 2022 Exam

Order of specific lectures may
change! Always check TUWEL!

Exams

First exam

- 13 June 2022, 12:00 – 14:00

Further exams

- October 2022
- December 2022
- March 2023

PROTIP

- You need to be able to explain 'things', e.g. what they are, why they are needed, what are pros and cons, etc.
- If you just memorise some acronyms then it won't help you much.

Lecture Materials

There is no single book to follow for this lecture!

Slides

- Meant for presenting!
- Contain links to papers/reports/webinars etc.
 - Check them!
- Take notes

DO NOT DEPEND ONLY ON SLIDES WHEN STUDYING FOR THE EXAM!

Lecture Materials

Literature (very basic)

- Barend Mons, Data Stewardship for Open Science: Implementing FAIR Principles, 10.1201/9781315380711
(available online through TU Bibliothek)
- Managing and sharing data by UK Data Archive:
<http://www.data-archive.ac.uk/media/2894/managingsharing.pdf>

Many links provided in the individual lectures – check them!

- <https://www.go-fair.org/fair-principles/>
- <https://www.fairsfair.eu>

Guidance for researchers from centres like the one at TUW:

- <https://www.tuwien.at/en/research/rti-support/research-data/center-for-rdm>

Overview of Exercises

194.045 Data Stewardship (UE 2,0) 2022S

- Practical exercises to get hands-on experience
- Programming skills essential

Exercise 1 (33%)

- Will be published this week in TUWEL
- Deadline ~ end of April

Exercise 2 (67%)

- Deadline ~ end of June
- Groups of 2

Submissions via TUWEL

To pass

- Min 35% of points on each exercise
- Min 50% of point overall to be positive

Late submission by 1 week: 10% points deduction

Overview of Exercises

Just to give an idea of what we did last year! Topics will be different this year!

Exercise 1

- use case - you pick an already existing experiment of yours or create one for this exercise;
- data management - you create a data management plan (DMP) and publish your experiment;
- machine-actionable DMP – you express your DMP using RDA DMP Common Standard.

Exercise 2

- Topic #1 SHACL constraints for DCSO
- Topic #2 maDMPs to support reviewers
- Topic #3 RO-Crates and Excel
- Topic #4 Meta-analysis of reproducibility of papers
- Topic #5 Analysis of DMPs from Exercise 1
- Topic #6 Core Trust Seal
- Topic #7 Repository content migration
- Topic #8 GitHub and Sustainability Metrics
- Topic #9 Python client library for TU Data

Selected Data Stewardship Activities and Collaborations

We teach what we learned from research projects and industrial collaborations

- e-Infrastructures Austria
- ROMOR: Research Output Management
- Timbus: Timeless Business Services (EU FP7 IP, via SBA)
- 4C: Collaboration to Clarify the Costs of Curation (FP7 CA, via SBA)
- APARSEN: (EU FP7 NoE, via SBA)
- Scape: Scalable Preservation Systems (EU FP7 IP)
- PLANETS (EU FP6 IP)
- DPE: Digital Preservation Europe (EU FP6 CA)
- ...

Center for RDM @ TUW

Center for Research Data Management

- FAIR Data Austria
 - [TU Data repository](#)
 - [TU Database repository](#)
 - [Tools for maDMPs](#)
 - [FAIR Office Austria](#)
- Research Data Management (RDM) Policy
- Support to all researchers on RDM

Research Data Management at TU Wien

You want to manage your data and your code according to the FAIR principles? You want to store, share and publish your data and need a repository? Your funding agency requires a data management plan? Use our services for research data management (RDM).

		
Technical RDM services und tools TU Data Repository, TU DMP Tool, TU gitLab	Center for RDM Training, consulting, projects	RDM infos and tips DMP handbook, RDM policy, funders' guidelines, basics



**FAIR DATA
AUSTRIA**

December 1, 2021 | Version 1.0.0

CLEF-IP 2013

Pirol, Florina

CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property

The CLEF-IP track ran from 2009 to 2013 and aimed to investigate IR techniques for patent retrieval. The track utilizes a collection of more than 1.3M patent documents (~2.6 million files) derived from EPO (European Patent Office) sources and EuroPCT Applications (more than 400K documents) published by WIPO (World Intellectual Property Organization). The collection contains documents in English, French and German with at least 150,000 documents in each language, all published before 2001.

There was one task in 2013. The first one was to find patent documents that are candidates to constitute prior art for a given claim taken from a patent document.

Files

- Document Collection**
The corpus consists of two parts. The first one is a set of XML files representing a total of over 1.3 million patent documents - this collection is to be used for the first task.
NOTE: the document collection is the same as the one published for CLEF-IP 2011, excluding images.
- Topics and Answers**
Both the training and the test topic sets contain also the relevance assessments for the topics.

Licenses
Creative Commons Attribution Non Commercial Share Alike 3.0 Unported
[Read more](#)

Files (14.3 GB)

Name	Size
01_document_collection.tgz	14.2 GB
02_topics.tgz	16.5 MB

Details

Licenses	cc-by-nc-sa-3.0	Creative Commons Attribution Non Commercial Share Alike 3.0 Unported
Resource type	Dataset	
Publisher	TU Wien	
Languages	English, French, German	
Alternate identifiers	DOI	10.48436/nw2xc-41j75
Related works	Is described by	10.1007/978-3-642-40802-1_25 (DOI) Publication
Upload information	Created: November 30, 2021 Modified: November 30, 2021	

Title of the dataset

Description of the dataset

Files for download

Persistent identifier (DOI):
[10.48436/nw2xc-41j75](https://doi.org/10.48436/nw2xc-41j75)

License: CC-BY-NC-SA-3.0

<https://researchdata.tuwien.ac.at>



FAIR
OFFICE AUSTRIA

We connect stakeholders from research communities and service providers. Together, we help to advance the FAIR principles.

About Contact




FAIR
OFFICE AUSTRIA

About Contact Deutsch

Information for Researchers

 Support Near You	 Best Practice Examples
 FAIR in the Research Proposal and Data Management Plan	 Infrastructures and Tools
 GO FAIR Implementation Networks	 Learn More About FAIR
 FAIRification Process	 Learn more about Research Data Management



FAIR
OFFICE AUSTRIA

About Contact Deutsch

Information for Service Providers

For us, **Service Providers** entail all entities who are responsible for developing FAIR services and tools, implementing FAIR principles, establishing FAIR practices and providing training on FAIR workflows and tools.

+ Service provider examples

 FAIR and the Research Data Alliance (RDA)	 GO FAIR Implementation Networks
 FAIR and the European Open Science Cloud (EOSC)	

<https://fair-office.at/>

GO-FAIR



GO FAIR

FAIR Principles Implementation Networks News Events Resources About GO FAIR

GO FAIR Austria office

Home > GO FAIR Initiative > GO FAIR Offices > GO FAIR Austria office

The **FAIR Office Austria** is a consortium of the three universities **TU Wien**, **Graz University of Technology** and the **University of Vienna**. The consortium members are all involved in the project FAIR Data Austria which is dedicated to implementing repositories and machine actionable data management plan tools at institutional level. The project also aims at establishing trained data stewards within Research Performing Organizations, thus supporting researchers in their process of integrating the FAIR principles into their activities. All consortium members are involved in global and national nodes of RDA and the European Open Science Cloud (EOSC). FAIR Office Austria was established in June 2021, see the **news item** announcing the office.

Address
TU Wien, **Center for Research Data Management**, Favoritenstrasse 13, 1040 Vienna, Austria
Graz University of Technology, **RDM Team**, Brockmanngasse 84, 8010 Graz, Austria
University of Vienna, **University Library**, Universitätsring 1, 1010 Vienna, Austria

Office Team
Barbara Sánchez-Solis (**contact**), Head of Center for Research Data Management
Ilire Hasani-Mavriqi (**contact**), Head of Research Data Management Team at TU Graz
Sarah Stryeck (**contact**), Data Steward at TU Graz
Alexander Gruber (**contact**), Data Steward at TU Graz
Susanne Blumesberger (**contact**), Head of Department Repositorymanagement PHAIDRA-Services
Tereza Kalová (**contact**), Project Manager of Department Repositorymanagement PHAIDRA-Services

Website
<https://fair-office.at>

<https://www.go-fair.org/go-fair-initiative/go-fair-offices/go-fair-austria-office/>

Research Data Alliance



Removing barriers in data sharing

Open community

- 9,600 members from 137 countries

Working and Interest Groups

- DMP Common Standards WG
- Dynamic Data Citation WG
- ...

RDA Recommendations

<https://www.rd-alliance.org>

An infographic titled 'Scalable dynamic-data Citation Methodology' is enclosed in a green border with decorative corner symbols (heart, club, diamond, spade). At the top center is the RDA logo. The main title is in green. Below it, the text 'The Challenge:' is followed by a paragraph: 'Supporting accurate citation of data subjected to change, for the efficient processing of data and linking from publications.' Then, 'What is the solution?' is followed by another paragraph: 'The Dynamic Data Citation Working Group developed a simple, scalable mechanism that allows the precise, machine-actionable identification of arbitrary sub selections of data at a given point in time irrespective of any subsequent addition, deletion or modification.' At the bottom center is an illustration of a blue calculator with a yellow data card on top. Below the illustration, it says 'Produced by: Data Citation (WGDC) WG' and provides the URL 'https://rd-alliance.org/groups/data-citation-wg.html'.

Innovationslehrgang Data Science und Deep Learning

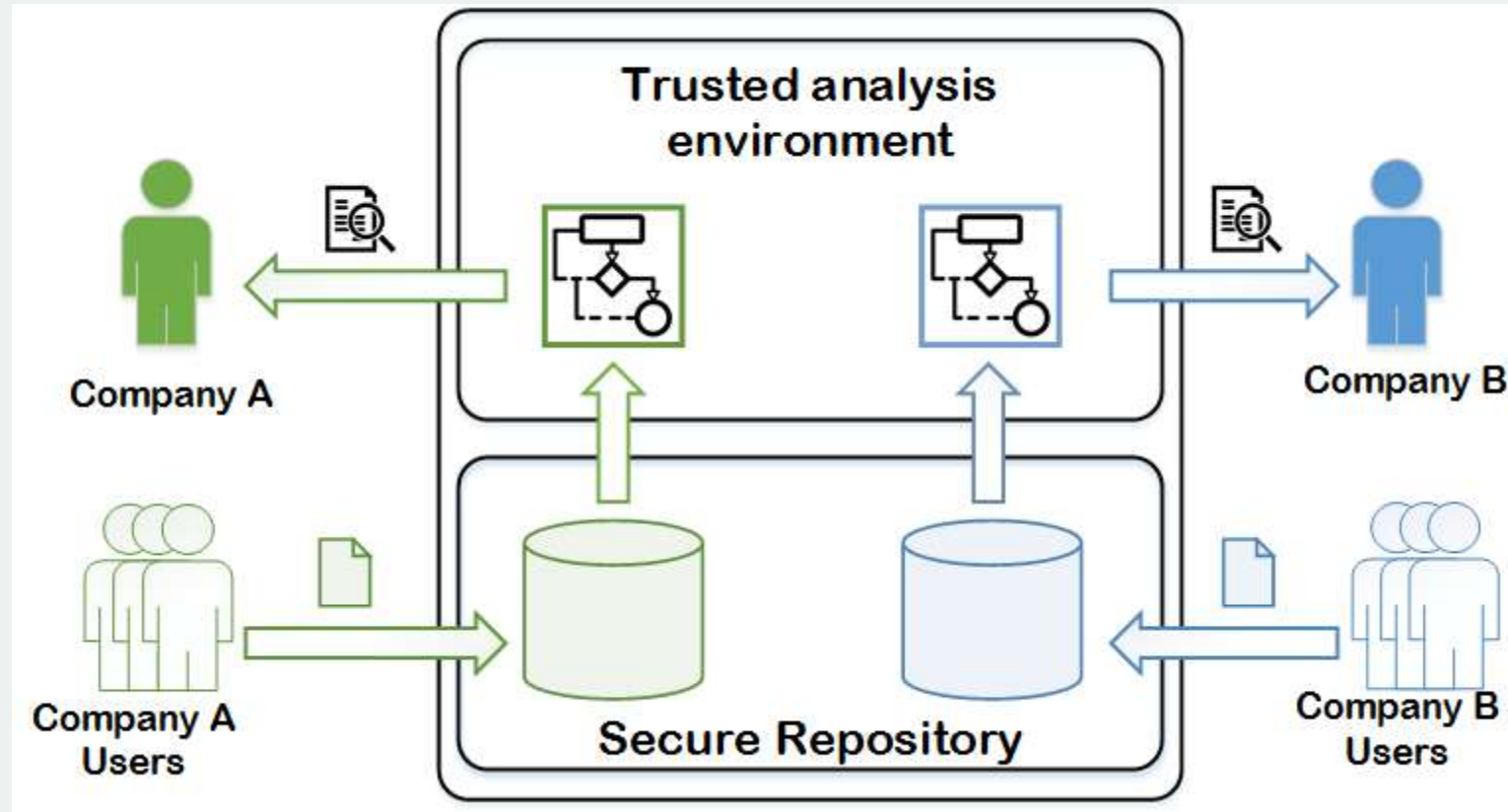
Data Stewardship was a module

<http://www.ifs.tuwien.ac.at/idsdl/>

iSDSL



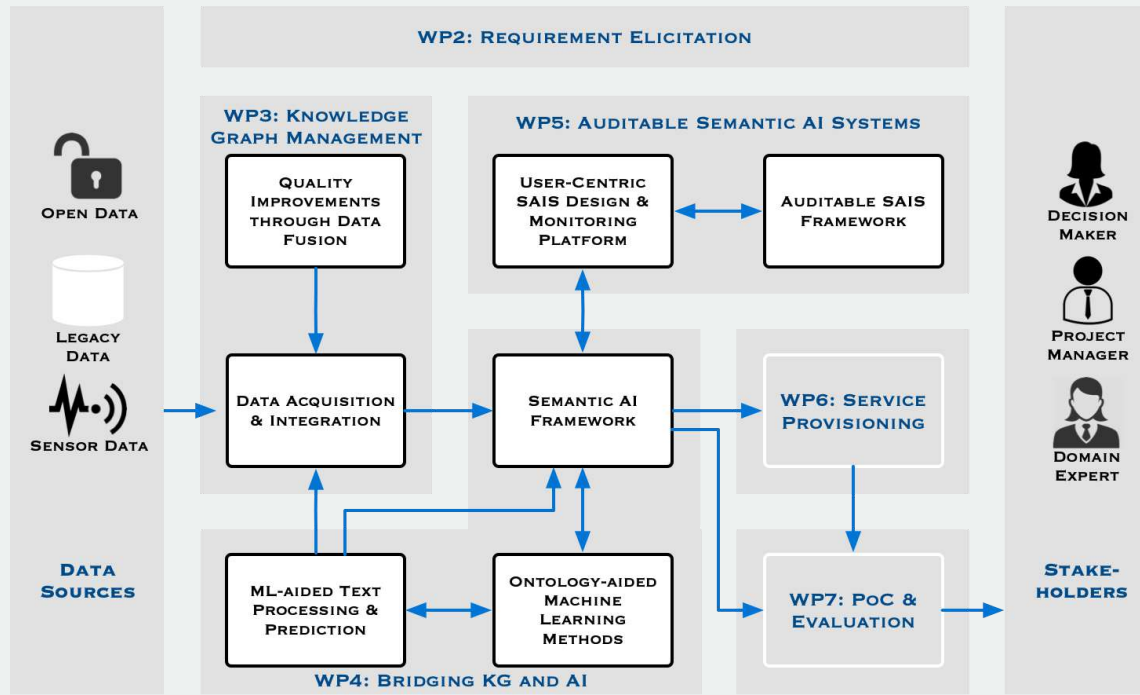
FFG WellFort



FFG OBARIS

Auditable Semantic AI Systems

Use case: *Predictive analytics of water quality*



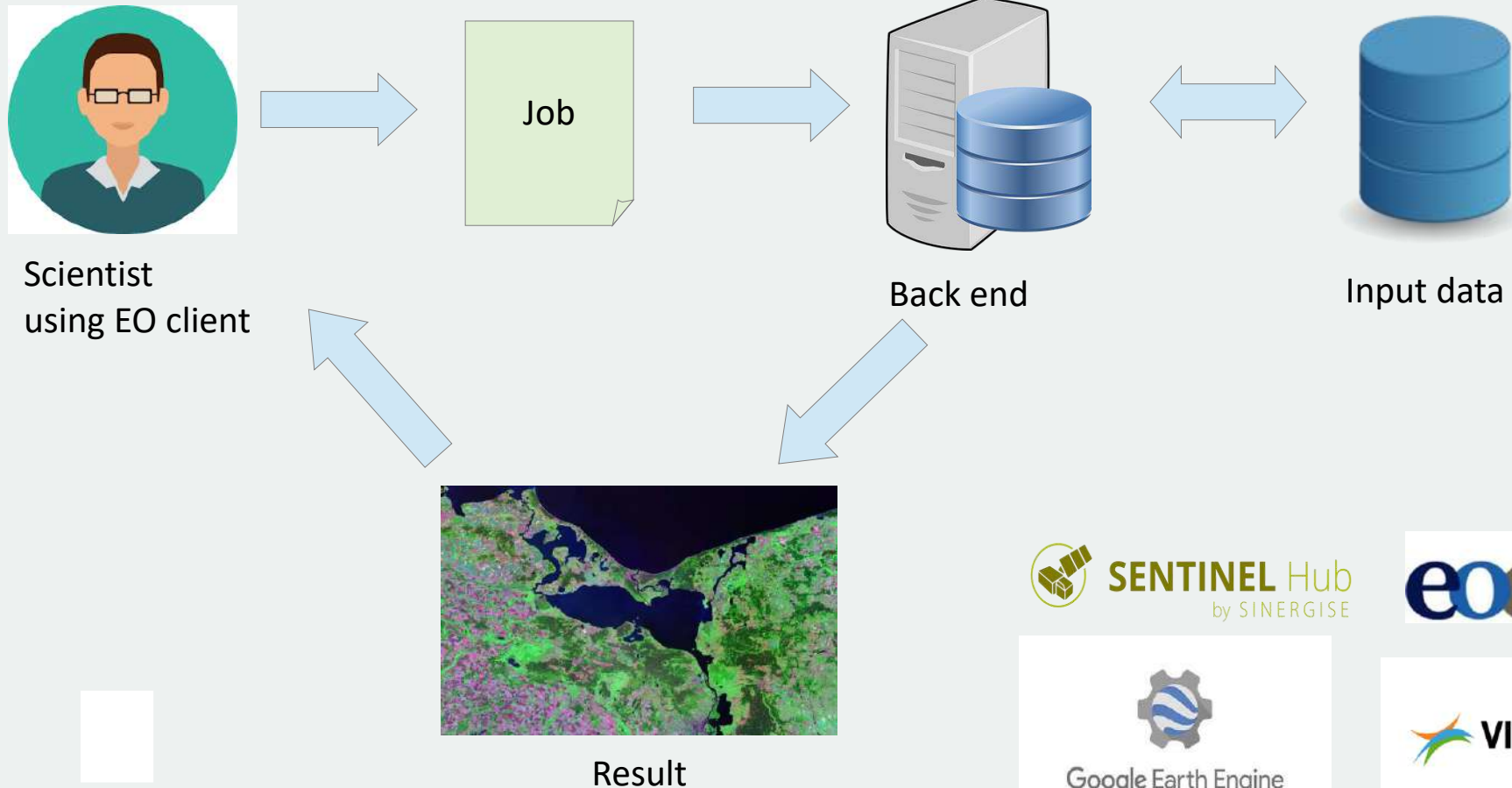
umweltbundesamt^U
PERSPEKTIVEN FÜR UMWELT & GESELLSCHAFT

SEMANTIC WEB COMPANY
linking data to knowledge

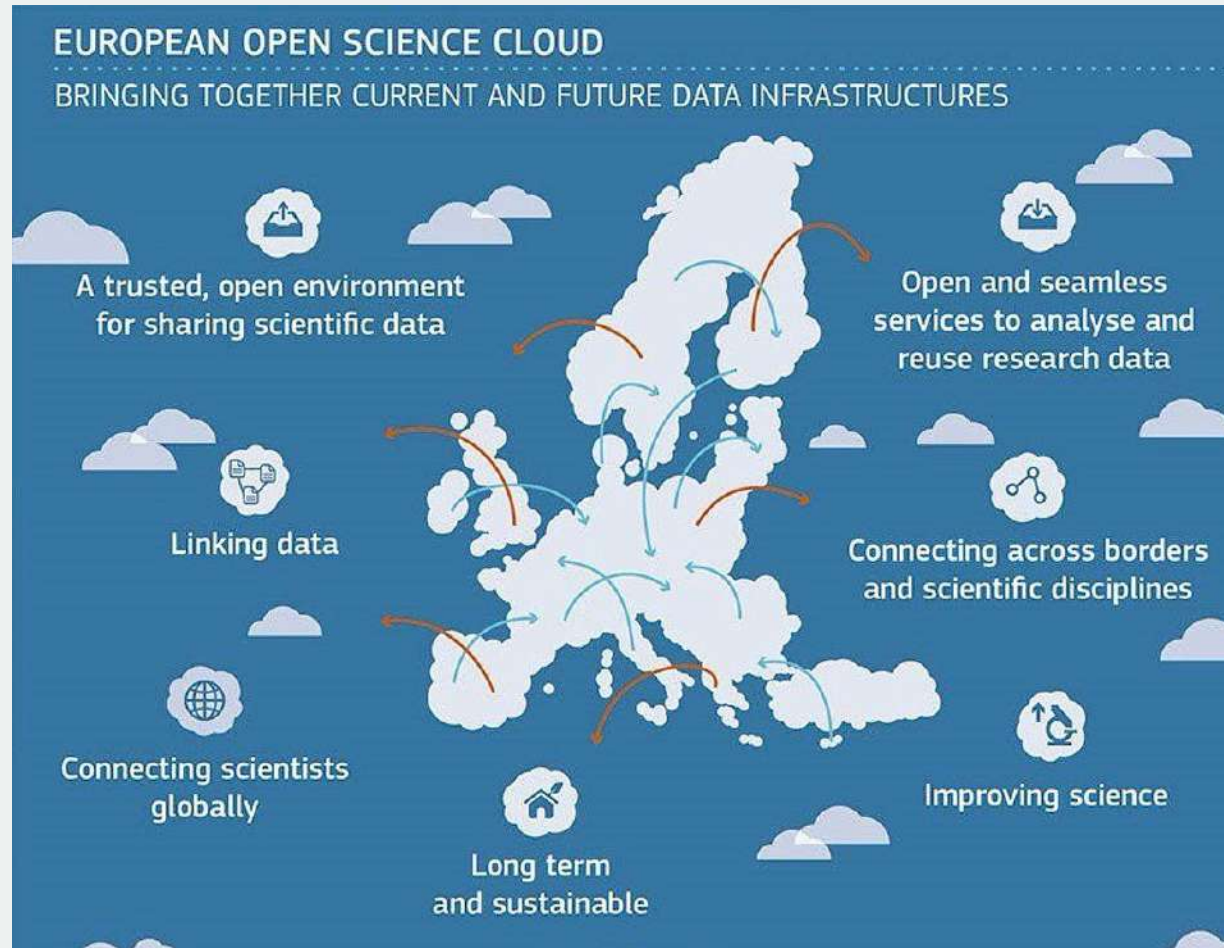
hpc
CONSULTING

<http://www.obaris.org>

EU openEO



European Open Science Cloud (EOSC)



EOSCsecretariat.eu

Setup and management of the EOSC Secretariat supporting the EOSC Governance

Funding opportunity



FEMTech Funding opportunity

- Funding by FFG (Austrian Research Funding Agency)
- Target group: **female students**
- Duration: 2-6 months, **paid**
- Can be used to work on a thesis or project (Praktikum)

In collaboration with SBA Research

- <https://www.sba-research.org>

More info

- <https://www.ffg.at/femtech-praktika>

First come, first served! Application has just opened!

Interested?

- ➔ Contact me (tmiksa@sba-research.org)



If you're looking for a topic...

We offer

- Master thesis topics
 - [“success stories”](#)
- Interdisciplinary projects, e.g. with other faculties

Focus

- Data management, versioning, reproducibility, auditability...

Various application domains

No synthetic problems

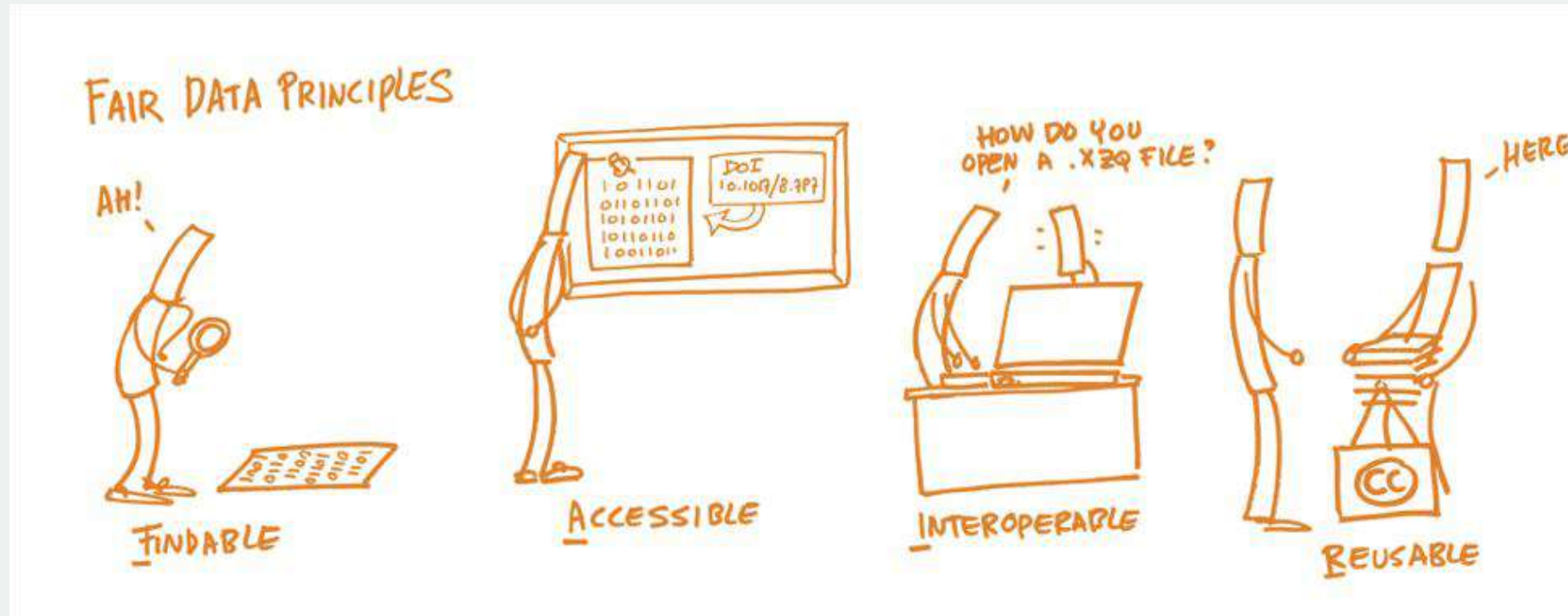
Usually the topics are not offered in TISS.

Just come and talk to us!

FAIR Principles

Introduction

FAIR Principles (very simplified :))



FAIR vs fair

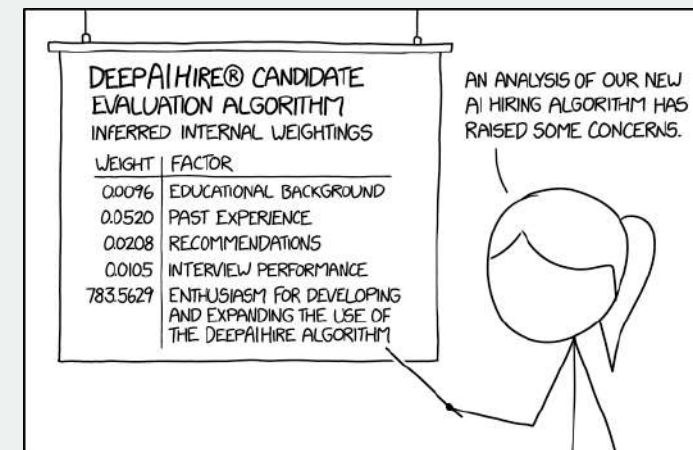
FAIR principles \neq Algorithmic fairness

To be FAIR

- To apply/use FAIR principles
- Focus on how data is managed, etc.

To be fair

- Evade bias
- Focus on design and implementation



<https://xkcd.com/2237/>

FAIR Principles



Home > FAIR Principles

> **FAIR Principles**

- > **F1: (Meta) data are assigned globally unique and persistent identifiers**
- > **F2: Data are described with rich metadata**
- > **F3: Metadata clearly and explicitly include the identifier of the data they describe**
- > **F4: (Meta)data are registered or indexed in a searchable resource**
- > **A1: (Meta)data are retrievable by their identifier using a standardised communication protocol**
- > **A1.1: The protocol is open, free and universally implementable**
- > **A1.2: The protocol allows for an authentication and authorisation where necessary**
- > **A2: Metadata should be**

In 2016, the **FAIR Guiding Principles for scientific data management and stewardship** were published in *Scientific Data*. The authors intended to provide guidelines to improve the **Findability, Accessibility, Interoperability, and Reuse** of digital assets. The principles emphasise machine-actionability (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention) because humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data.

A practical "how to" guidance to go FAIR can be found in the **Three-point FAIRification Framework**.

Findable
The first step in (re)using data is to find them. Metadata and data should be easy to find for both humans and computers. Machine-readable metadata are essential for automatic discovery of datasets and services, so this is an essential component of the **FAIRification process**.

F1. (Meta)data are assigned a globally unique and persistent identifier

F2. Data are described with rich metadata (defined by R1 below)

F3. Metadata clearly and explicitly include the identifier of the data they describe

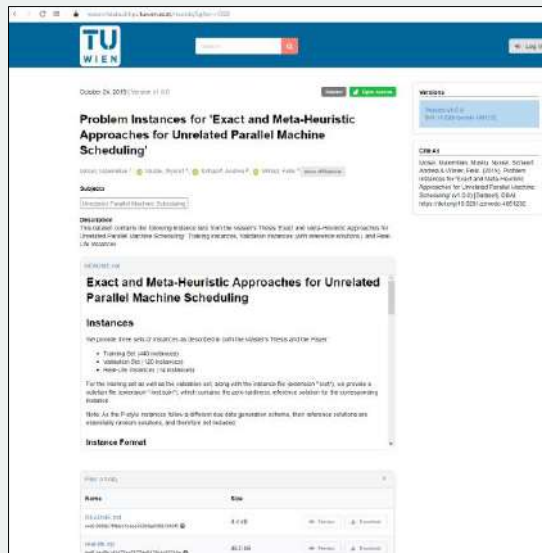
F4. (Meta)data are registered or indexed in a searchable resource

Accessible
Once the user finds the required data, she/he needs to know how can they be accessed, possibly including authentication and authorisation.

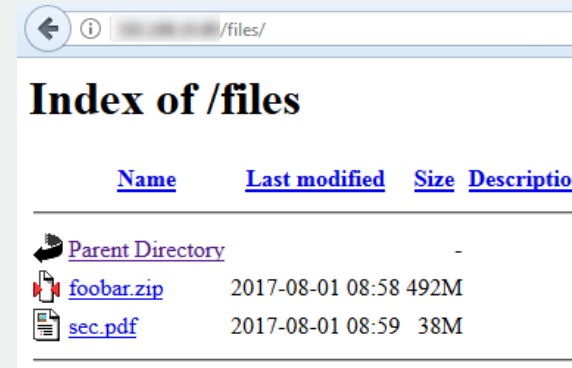
A1. (Meta)data are retrievable by their identifier using a standardised communications protocol

<https://www.go-fair.org/fair-principles/>

Findable – simplified examples



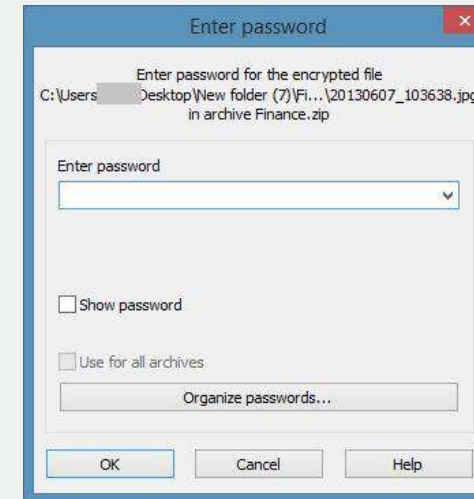
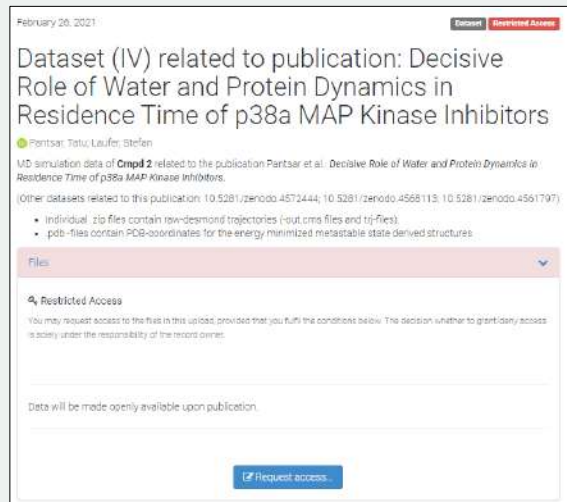
Data repository



Personal website



Accessible – simplified examples



Restricted access, but a clear way to request access

Interoperable – simplified examples

Yes

- XML following known XSD Schema
- MP3 for audio recordings
- Data model using common vocabularies



No

- Custom XML without any documentation
- M4P (Apple) for audio recordings
- Custom fields in data model with poor documentation



Reusable – simplified examples



Jährliche Personeneinkommen

Die knapp 4,6 Mio. unselbständig Erwerbstätigen (ohne Lehrlinge) erzielten 2019 ein mittlere Bruttojahreseinkommen von 29.498 Euro. Die Einkommen der Frauen streichen mit 22.805 Euro im Mittel nur 63,8% des Einkommens der Männer (35.841 Euro), wobei Frauen viel häufiger Teilzeittätige sind. Die mittlere Nettojahreseinkommen belaufen sich auf 22.894 Euro (Frauen: 19.233 Euro, Männer: 25.938 Euro).

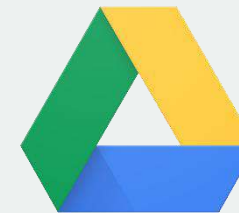
Warten die Einküsse von Teilzeit und nicht-ganzjähriger Beschäftigung ausgeglichen und nur Personen berücksichtigt, die laut Lohnverordnen Vollzeit beschäftigt sind und im Jahr 2019 mindestens 150 Tage im Jahr unselbständig erwerbstätig waren (ohne Lehrlinge), so beträgt das mittlere Bruttojahreseinkommen der Frauen 36.326 Euro, während Männer im Mittel 45.900 Euro verdienen. Der relative Einkommensanteil der Frauen am mittleren Einkommen der Männer stieg somit auf 85,7%.

Bei einer Unterteilung nach Beschäftigungsgruppen müssen wegen der Anteile an Teilzeittätigen weitere strukturelle Unterschiede berücksichtigt werden. Insbesondere ungleiche Anteile nicht-ganzjähriger Beschäftigung sowie unterschiedliche Qualifikations- und Altersstrukturen. Da sind beispielsweise Beamtinnen und Beamte deutlich älter, weisen ein höheres Ausbildungsniveau auf und sind kaum besetzungsabhängig. Bei einer Einteilung der Beschäftigungsgruppen nach dem Median der Bruttojahreseinkommen stellen männliche Beamte an der Spitze (56.772 Euro) gefolgt von Dozenten (56.200 Euro), männlichen Angestellten (47.373 Euro) und männlichen Vertragsbediensteten (42.112 Euro). Weibliche Vertragsbedienstete (32.362 Euro), männliche Arbeiter (21.454 Euro) und weibliche Angestellte (21.448 Euro) im Rahmen deutlich weniger, das Schwächste bilden die Arbeiterinnen (12.883 Euro). Betrachtet man die Beamtinnenverdienste (oberste Decke), so liegen die männlichen Angestellten mit 40.844 Euro knapp hinter den männlichen Beamten mit 48.742 Euro, deutlich darüber stehen die bestbezahlenden Beamtinnen mit 33.519 Euro.

↓ mehr

Tabellen	Graphiken	Dokumentationen	allgemeine Auskünfte
Brutto- und Nettojahreseinkommen der unselbständig Erwerbstätigen 1997 bis 2019			
Brutto- und Nettojahreseinkommen der ganzjährig Vollzeitbeschäftigten 2004 bis 2019			
Brutto- und Nettojahreseinkommen der Personalarbeiter und Personalarbeiterinnen 1997 bis 2019			
Bruttojahreseinkommen der unselbständig Erwerbstätigen 2019			
Nettojahreseinkommen der unselbständig Erwerbstätigen 2019			
Nettojahreseinkommen der ganzjährig Vollzeitbeschäftigten 2019			
Nettojahreseinkommen der ganzjährig Vollzeitbeschäftigten 2018			
Bruttojahreseinkommen von Frauen und Männern 2019			
Endlohnjahreseinkommen von Frauen und Männern nach Bundesländern 2010			
Brutto- und Nettojahreseinkommen nach Altersgruppen 2019			

Trusted source, permission to reuse, well defined meaning of terms used



Provenance and permissions not clear

Machine-actionability lies at the heart of FAIR!



Not machine-actionable

A screenshot of a PDF viewer window showing a table with the following content:

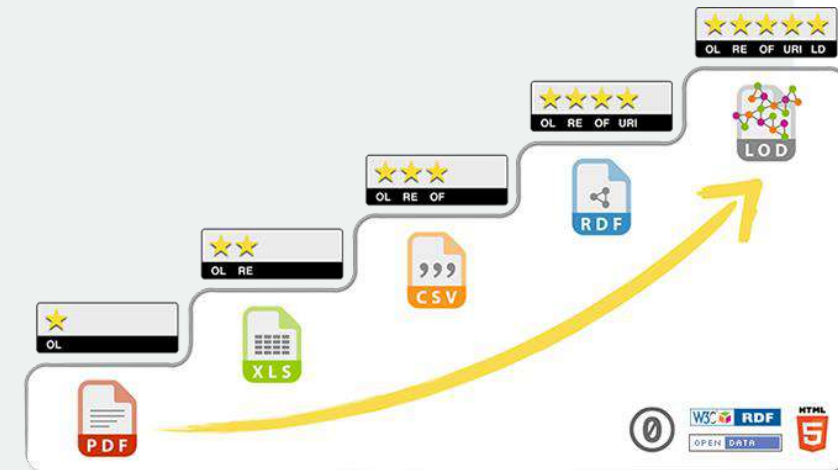
Temperature forecast for Galway, Ireland	
Day	Lowest Temperature (°C)
Saturday, 13 November 2010	2
Sunday, 14 November 2010	4
Monday, 15 November 2010	7



Machine-actionable

A screenshot of a spreadsheet application showing the same data as the PDF, but in a machine-actionable format:

	A	B
1	Temperature forecast for Galway, Ireland	
2		
3	Day	Lowest Temperature (°C)
4	Saturday, 13 November 2010	2
5	Sunday, 14 November 2010	4
6	Monday, 15 November 2010	7
7		



<https://5stardata.info/en/>

FAIR

More in a dedicated lecture!

In the meantime you can watch:

- Let's Make Our Data FAIR! Webinar for GO-FAIR

<https://www.youtube.com/watch?v=dEV2Hnraqal>

Next lecture

14 March 2022

Digital Preservation

- How to keep data available for the long term?





Digital Preservation Introduction

Andreas Rauber

Department of Software Technology and
Interactive Systems

Vienna University of Technology

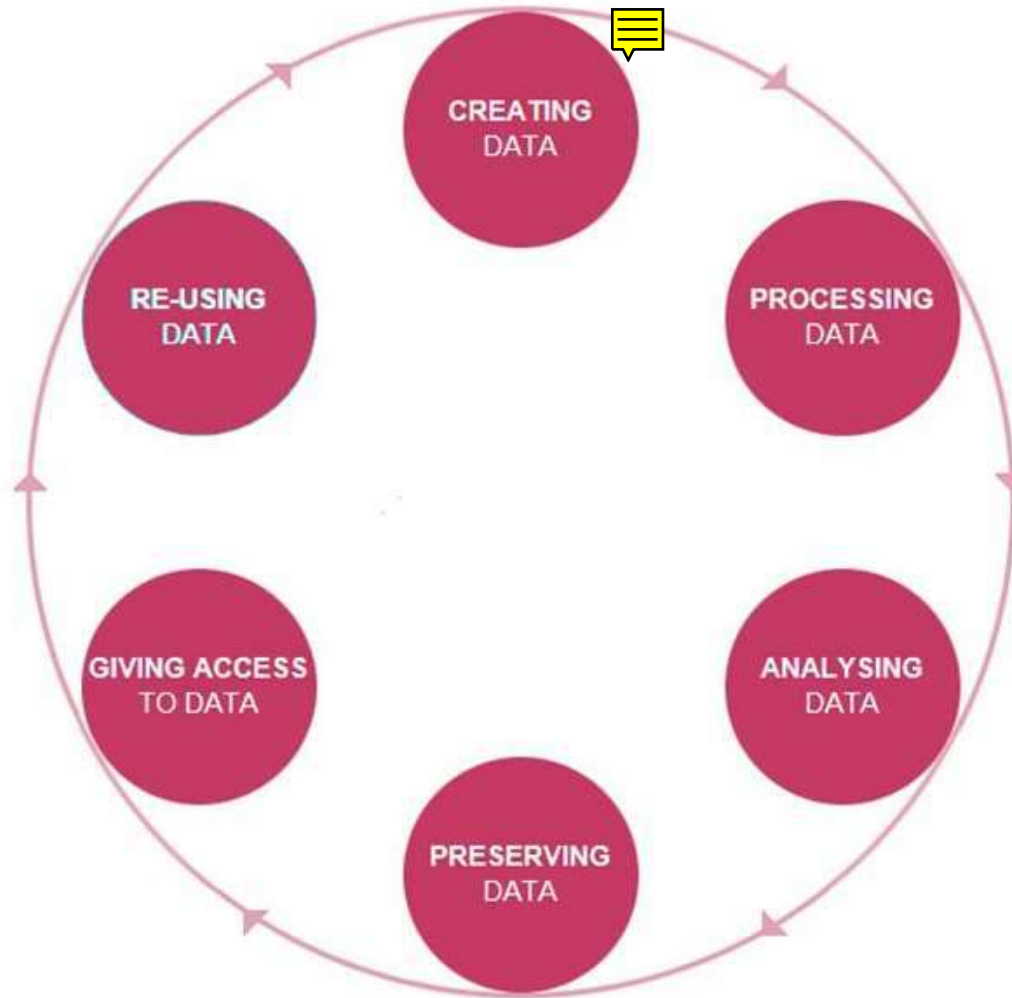
<http://www.ifs.tuwien.ac.at/~andi>

.....



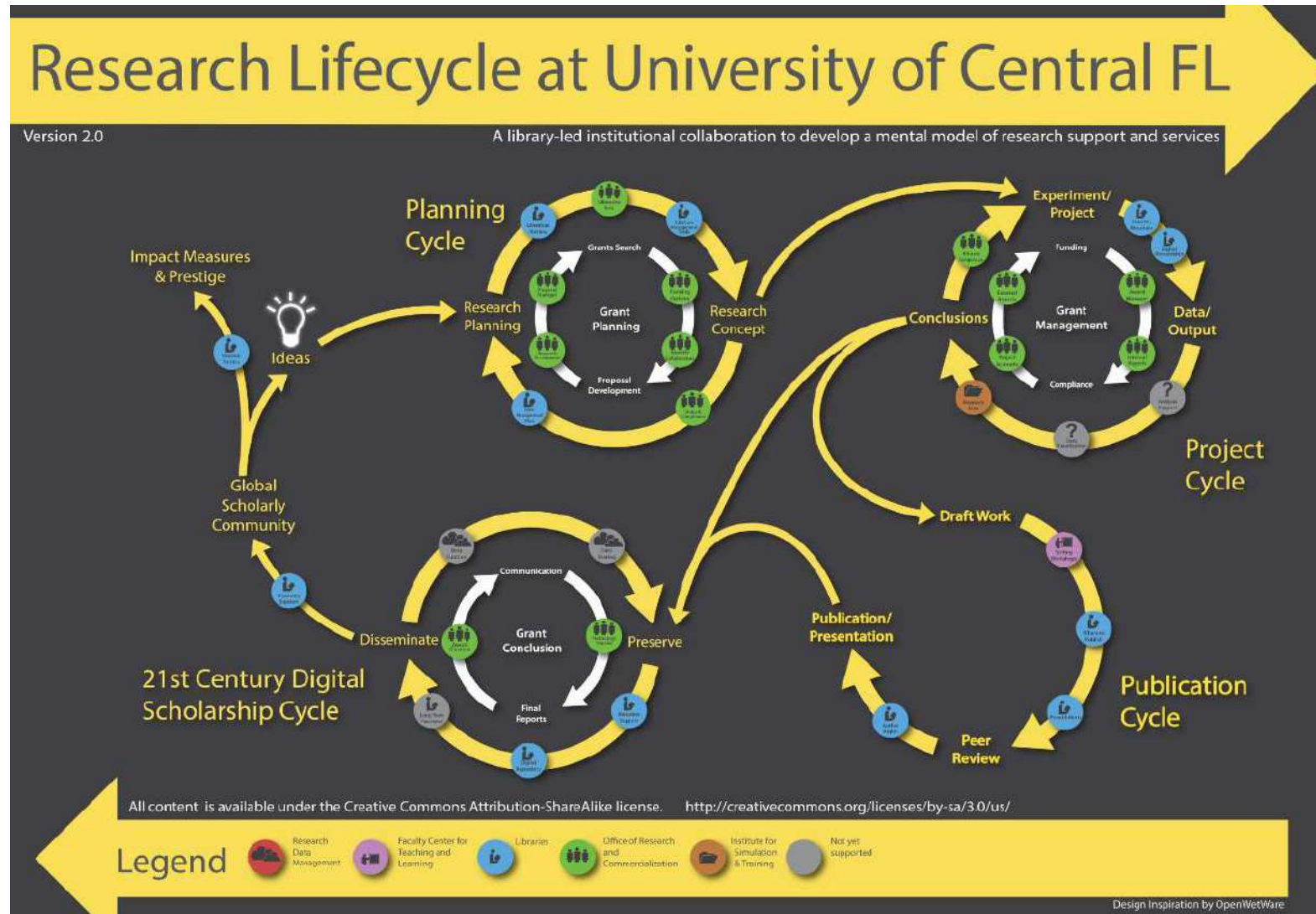
FACULTY OF **INFORMATICS** 1

UK Data Archive Lifecycle model



.....

University of Central Florida Lifecycle Model



Overview

-
- What are the challenges in Digital Preservation?
 - What methods can we use to counter them?
 - What's the issue with file formats?
 - What do we have in our repository?
-



Why do we need Digital Preservation?

Questions / discussion:

- What is *Digital Preservation*?

.....

Why do we need Digital Preservation?



.....



Why do we need Digital Preservation?

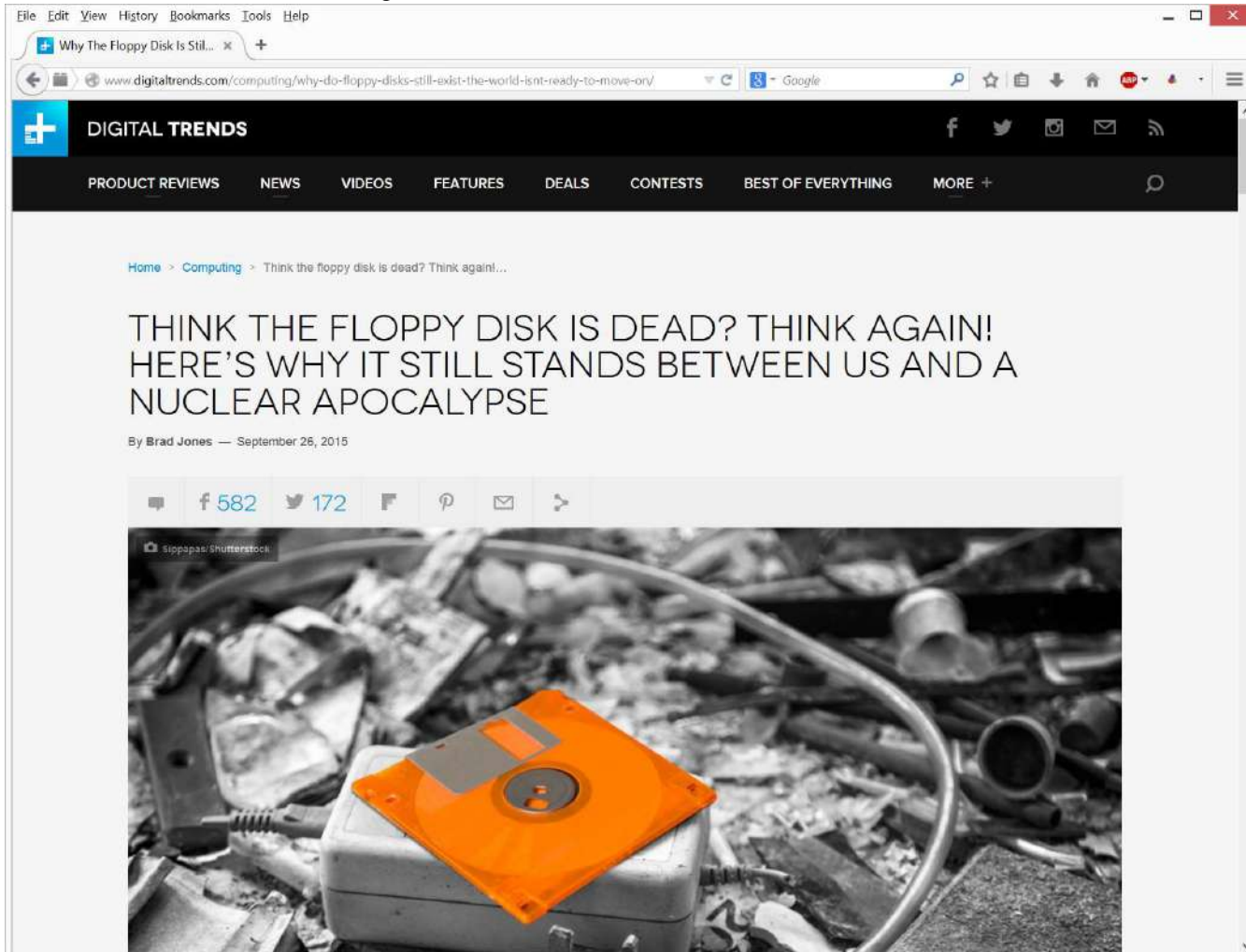
1. Physical Preservation (Bit-stream preservation)

- Transferring to current storage systems
 - note: transfer may not be trivial (file systems, encodings, relative references, copy protection,...)
- Ensure redundancy
 - technologically
 - geographic spread
- Access, security
- Error detection, recovery, disaster planning

.....

Why do we need Digital Preservation?

Just as a curiosity:

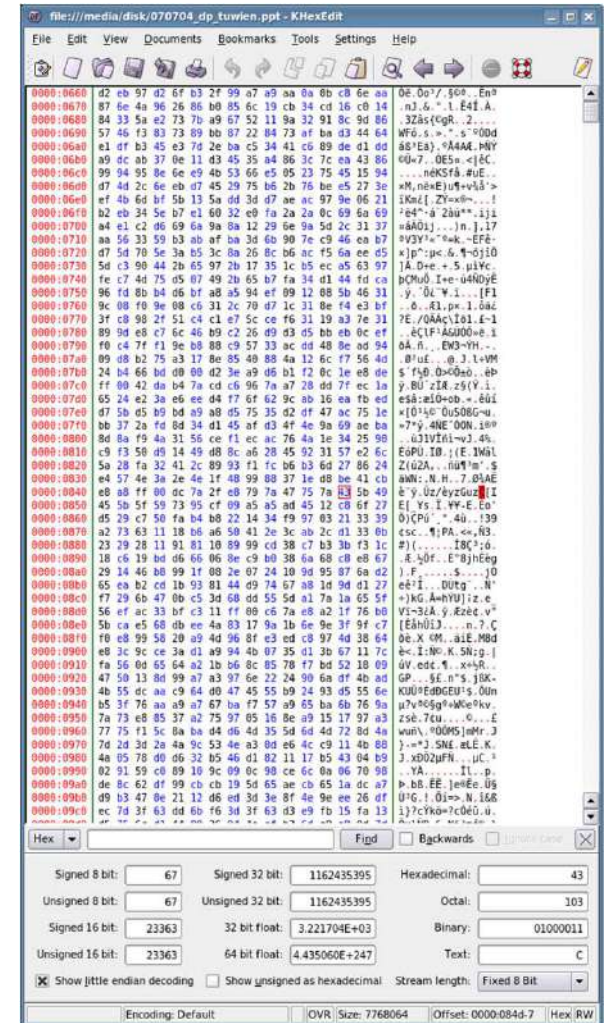
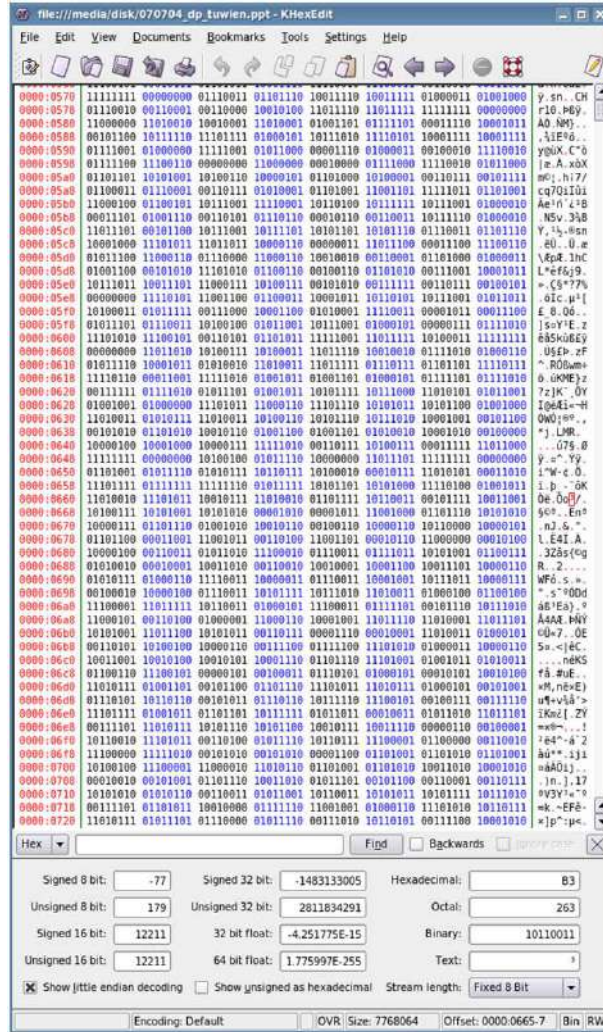
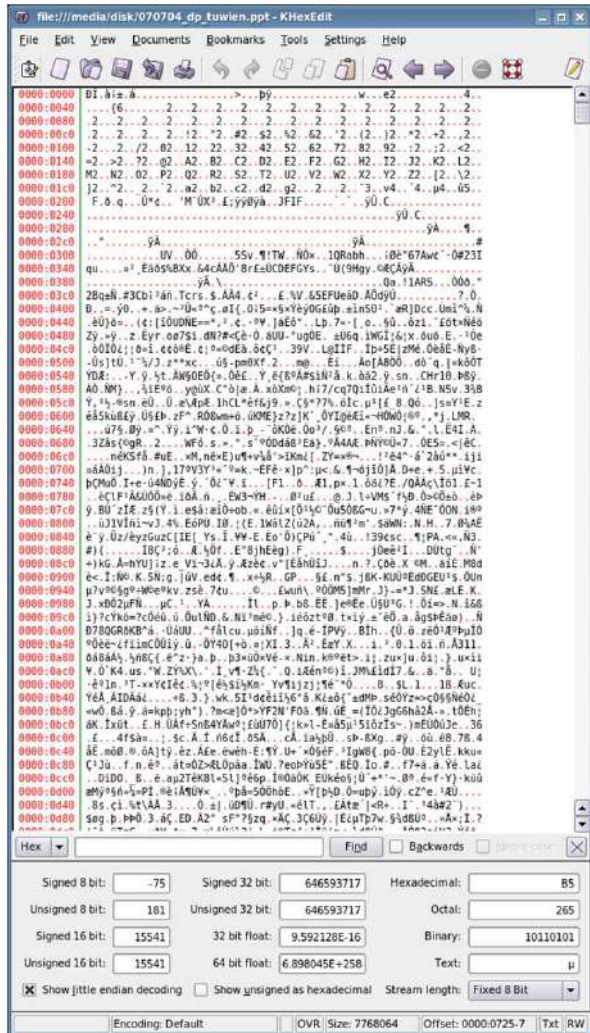


<http://www.digitaltrends.com/computing/why-do-floppy-disks-still-exist-the-world-isnt-ready-to-move-on/>

.....



Why do we need Digital Preservation?





Why do we need Digital Preservation?

2. Logical Preservation

- Digital Objects require specific environment to be accessible :
 - Files need specific programs
 - Programs need specific operating systems (-versions)
 - Operating systems need specific hardware components
- SW/HW environment is not stable:
 - Files cannot be opened anymore
 - Embedded objects are no longer accessible/linked
 - Programs won't run
 - Information in digital form is lost (usually total loss, no degradation)
- Digital Preservation aims at maintaining digital objects authentically usable and accessible for long time periods.

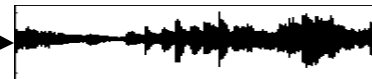
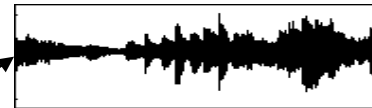
.....

Strategies for Logical Preservation

Another way of viewing this...



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



.....



https://en.wikipedia.org/wiki/Konzerthaus,_Vienna



https://en.wikipedia.org/wiki/Odeon_of_Herodes_Atticus



https://commons.wikimedia.org/wiki/File:Stereoanlage_Vision_2000.jpg

.....



Why do we need Digital Preservation?



Homann Heirs Map
Wikimedia 1747

Cary Map of Au
Wikimedia 1801

Austria 1999 CIA map
Wikimedia, 1999

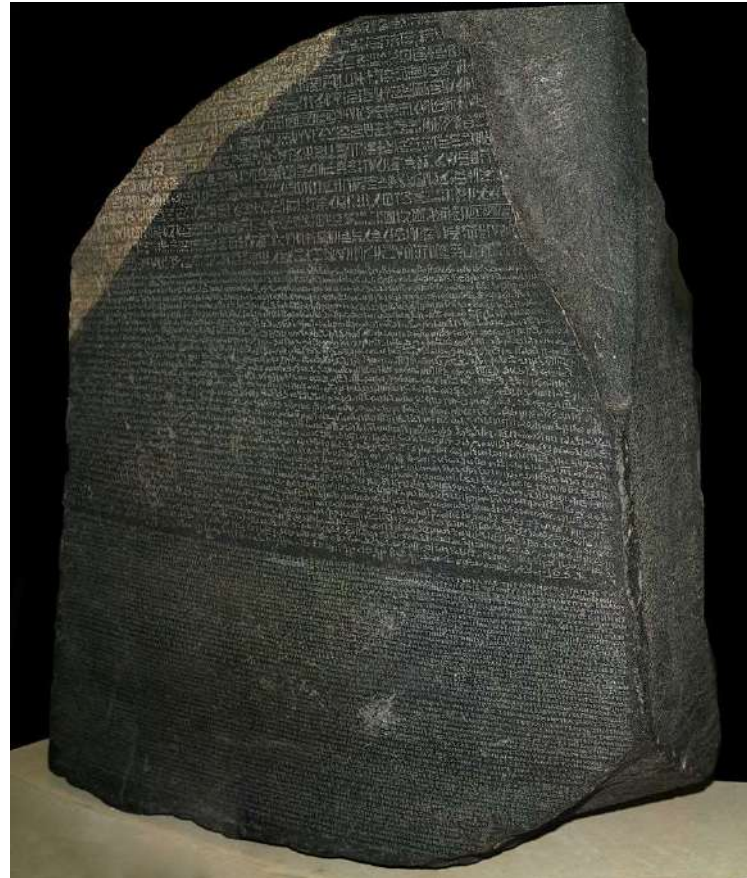
Mitchell Map of Austria, Hungary and Transylvania
Wikimedia 1850

3. Semantic Layer: information object

- How to interpret the data (information?) in the objects?
 - terminology changes:
changes in country names, borders, connotation of words,...
 - concept changes:
drunk driving: before 1998: 0.8‰ , afterwards 0.5‰
 - transformations: currencies/exchange rates, sensor resolutions,
 - provenance: actions applied to objects
sources: who? / which sensor?, transformations, post-processing
 - context of objects:
understanding the context of decisions, side-effects, quotations,
calibration timestamps
- For preserving digital information, all 3 layers need to be addressed

Why do we need Digital Preservation?

One of the most famous examples...



https://commons.wikimedia.org/wiki/File:Rosetta_Stone.JPG

.....

Why do we need Digital Preservation

- The goal of Digital Preservation is to **maintain digital objects accessible and usable in an authentic manner for a long term** into the future.

Why do we need Digital Preservation?

Questions / discussion:

- What is *digital data*?
- What is *digital storage*?
- What do we mean by
 - *accessible*?
 - *authentic*?
 - *long-term*?

.....

Why do we need Digital Preservation?

- Essential for all digital objects
 - Office documents, accounting, emails, ...
 - Scientific datasets, sensor data, metadata, ...
 - Applications, simulations, business processes, ...

- All application domains
 - Cultural heritage data
 - eGovernment, public administration
 - Science / Research
 - Industry
 - Health, pharmaceutical industry
 - Aviation, control systems, construction, ...
 - Private data
 - ...

.....

Why do we need Digital Preservation

3 levels of threat / preservation

1. Bit rot – physical preservation / bit preservation
Physical Layer: how to keep the 0's and 1's
2. Object formats – logical preservation
Logical Layer: how to remain able to open a file, run a program
3. Authenticity, interpretability – semantic preservation
Semantic Layer: how to ensure we can understand/interpret data correctly

- What can we do?

.....

Overview

-
- What are the challenges in Digital Preservation? (recap)
 - What methods can we use to counter them?
 - What's the issue with file formats?
 - What do we have in our repository?
-

Bit-level preservation

- Maintain bit-sequence
- Redundant storage:
 - Lockss: lots of copies keeps stuff safe
 - Cloud
- Distributed storage – physically separated
- Different technologies / platforms / production batches
- Controlled storage conditions
- Regular maintenance: tape rewinding, disc spinning, ...
- Maintain devices for accessing storage!
- Trade-off capacity, energy, effort

.....

Bit-level preservation

Questions / discussion:

- How long do tapes / CDs / DVDs / HDDs / SSD last?
- What are the costs of bit-level preservation?
- What are the logistic challenges?
- Is a DVD that lasts for 200 years a solution?
- What would be the most durable storage technologies?
- What is "digital storage"?
- Distribution and Trust?
- Are we allowed to store redundantly? in the cloud?
 - Copyright
 - Copy protection
 - Distributed objects, referenced via URL? DOI?

Why do we need Digital Preservation

3 levels of threat / preservation

1. Bit rot – physical preservation / bit preservation
Physical Layer: how to keep the 0's and 1's
2. Object formats – logical preservation
Logical Layer: how to remain able to open a file, run a program
3. Authenticity, interpretability – semantic preservation
Semantic Layer: how to ensure we can understand/interpret data correctly

- How? What can we do?


.....

Technology Museum

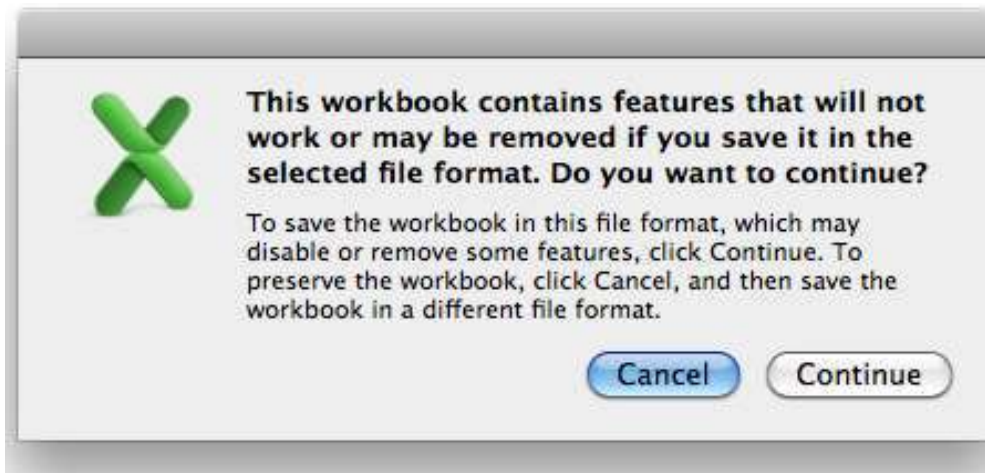
- Keep the hardware (drives, computer,...)
- + Maintains full functionality
- + Creates time buffer to develop more permanent strategies
- + Requires detailed documentation of HW and SW, but this also helps
- + Only strategy for some types of objects? (which?)
- Economically and technically infeasible to maintain spare parts forever
- Requires huge "museum"
- Requires highly specialized know-how for all platforms and software



Migration

- Transform into different format 
- Continually or on demand (Viewer)
- + Widely used
- + Possibility to compare at time of migration
- + Resulting objects are always accessible
- Possibly undesired changes during migration
- Needs to be repeated again and again

Strategies for Logical Preservation

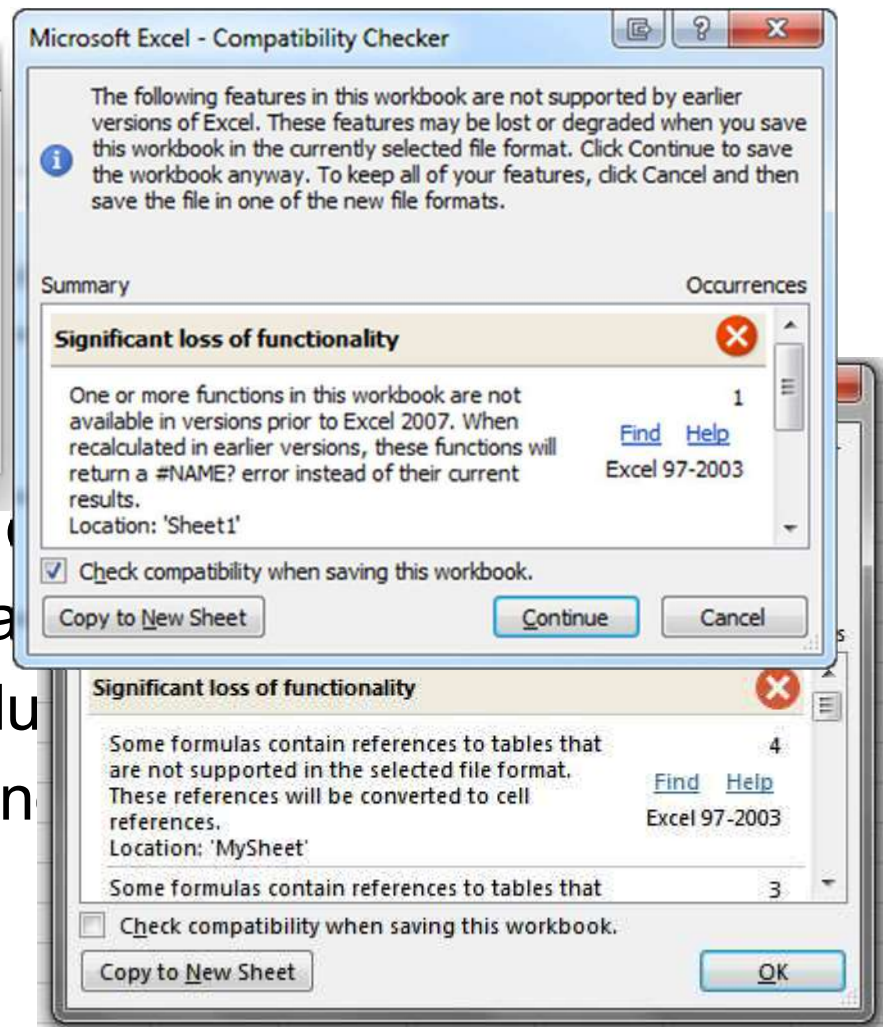
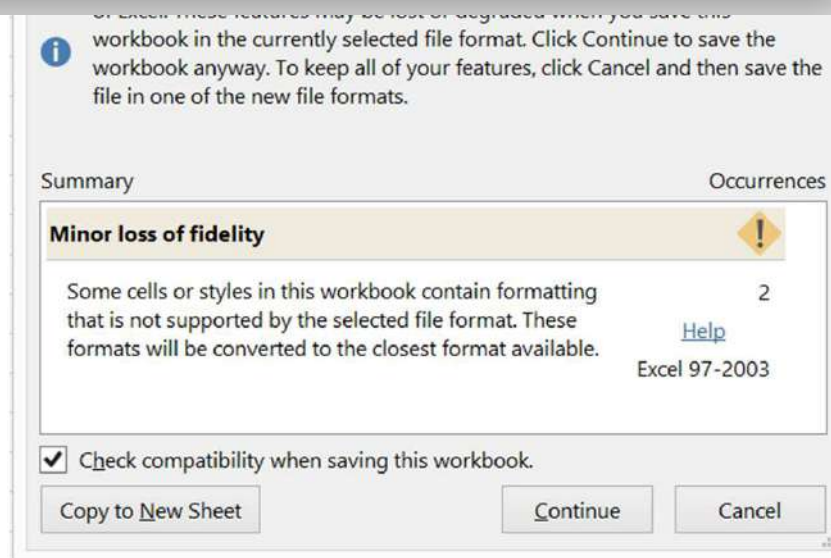


+

+

-



-



<https://support.office.com/en-us/article/Use-Office-Excel-2010-with-earlier-versions-of-Excel-2fd9ffcb-6fce-485b-85af-fecfd651a5ac>

.....

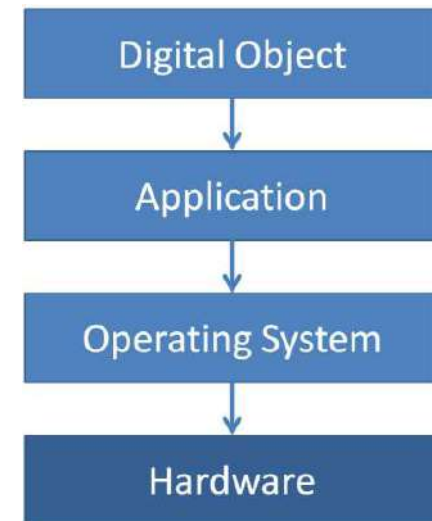
Emulation

- Emulation of Hardware or Software (OS, application)
- + Widely used principle
- + Many emulators available
- + Potentially preserving complete functionality
- + *Document is unchanged*
- *Document is unchanged* 
- Complex technology, lot of research required
- Requires detailed documentation of the system
- Requires experience how to interact with emulated historic system in the future
- Emulators must be migrated as well 
- Emulators potentially erroneous (Complexity)

.....

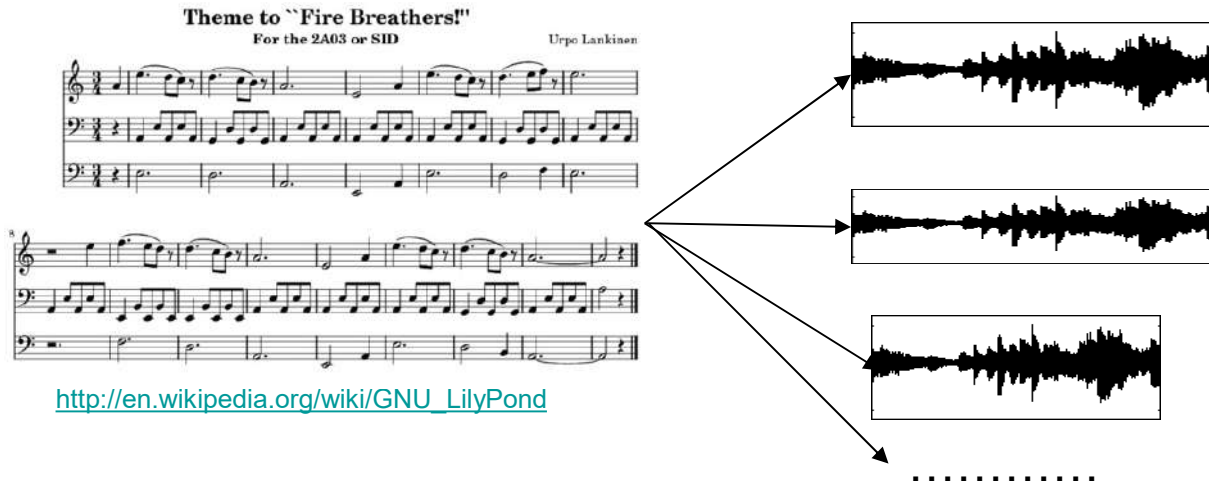
Excursion: Emulation vs. Migration

- Different on the pragmatic level, but conceptually identical
- Change occurs somewhere in the viewpath
- Have basically the same advantages/disadvantages and characteristics
- None of them guarantees identical rendering/performance of digital objects
- Many variants (e.g. viewer, virtualization)
- Need to be evaluated the same way



Strategies for Logical Preservation

Emulation vs. Migration – remember this?



https://en.wikipedia.org/wiki/Konzerthaus,_Vienna



https://en.wikipedia.org/wiki/Odeon_of_Herodes_Atticus



https://commons.wikimedia.org/wiki/File:Stereoanlage_Vision_2000.jpg


.....

Standardization

- Using open or de-facto standards
- + Simplifies DP process
- + Many tools available
- + Tools for standards are easier to build also in the future
- Significant effort required for standardization
- Loss at converting into standard
(who is responsible?)
- Some object types cannot be standardized

Strategies for Logical Preservation


Standardization - Excursion into file formats Proprietary vs. Open

- Proprietary
 - Documentation mostly not available
 - License and patent rules
 - License agreements subject to change
 - Restrictions for use and modifications may apply
- Open
 - Documentation available!
 - Unlimited use
 - No license fee
 - Open for modifications
 - No patent owners
- But: sometimes proprietary may be better than open - **why?** 
- Is the concept of "file formats" still useful?

Limiting Accepted Formats

- Similar to standardization
- + Reduces challenge to smaller number of formats
- Does not solve the problem
- Limits the type of objects that can be accepted
- Potential loss at conversion
- Requires strict control of formats (and what's in them!)

Data/Information Extraction

- Create abstract representation of information (e.g. databases or documents -> XML) 
- + Independent of specific infrastructure
- + Many tools available
- + Easier to develop tools in the future
- High effort to develop tools for specific abstraction scenario
- Limited functionality of tools designed to interpret information, many aspects not preservable
- Cannot be applied to all types of objects

Encapsulation

- Add metadata, software,... (representation information) to object („onion“)
- + Simplifies search for preservation solution on demand, offering several potential layers
- + Always allows for the application of several other strategies at different levels
- Does not solve the problem
- Even with all information encapsulated we may not be able to find a solution

Universal Computing Platform

- Example: UVC: Universal Virtual Computer (IBM)
- Abstract virtual machine, intermediate platform that can be implemented on many other platforms
- + Works for documents and software
- + A kind of standardization for platform, reduces development effort
- + Can test solution at time when being developed
- Pretty complex (cf. Java, but that's still simple)
- High effort at time of preservation
- Requires cooperation of the producers of information
- High risk of losing aspects of information

Backwards Compatibility and Version Migration

- current SW reads old versions and performs migration
- + Usually available
- + Creates time buffer for more permanent solutions
- + sometimes equal or better functionality
- Doubtful whether this will work for a long time (why?)
- Each change might lead to unwanted changes
- No guarantee from part of the producer of the SW

Strategies for Logical Preservation

Viewer



- Migration on demand, interpretation by Viewer software
- + Original datastream unchanged, interpreted directly
- + No continuous migration
- + No cumulative errors
- Viewer sometimes cannot process all (parts of) objects
- Increasing time delay when developing viewers
- Viewer SW must be carried along with technology changes
- Hard to evaluate whether viewer is correct

Non-digital Strategies

- Printing to paper, microfilm, ...
- + Requires transformation to readable form -> stable
- + Coding of digital data is possible
- + Lots of experience in handling analog data carriers
- + High stability -> Bit-stream Preservation
- Loosing functionality, loosing advantage of digital technology
- Not applicable for all objects
- High costs for preserving some of the analog data carrier material, low storage density, ...
- Even this can be “buggy” (Xerox bug, manipulation)



Data Recovery, Data Archeology

- Analysis of bit-stream to interpret data, digital forensics
- + Probably only approach to recover "lost" information
- No guarantee that it works
- Without sufficient documentation close to "guessing"
- Extremely high costs per object
- Hard to guess whether it may be successful for a given object

Summary


- Changing object, environment
- Loss upon migration / emulation
- Decision of what to preserve → **Significant Properties!**
- How to detect/document what you lost?
- Range of strategies available, none is perfect
- Combination of strategies
- No solution forever -> DP is a process!

Logical Preservation

- Preservation Planning
- Identify objects at risk
- Standardization reduces risk (why?)
- Apply preservation actions such as migration / emulation / HW-museum
- Identify what you need to preserve (significant properties)
- Identify suitability of tools
- Find out what you can preserve / what you loose
- Do it, document it, verify it, monitor it


Logical Preservation

Questions / Discussion:

- What are the problems of logical preservation?
- What is the optimal strategy? 
- What is the optimal strategy for a specific object?
- What is a good format / platform (e.g. to migrate to)?
- What are characteristics of good formats/platforms/... ?
- How can we identify objects at risk?
- When is a format "more/less risky"?
- What is a file format?
- How can we find out what we loose with a strategy?

Logical Preservation


Questions / Discussion (2):

- What is the difference between emulation and migration? Are they different? Are they not different?
- What are the significant properties of an object / process? 
- “I want to preserve everything” – (how) can we do this?
- What is the “original object”?
- What is the complexity of each strategy? Costs? Effort?
- What know-how do we need to decide on a strategy?
- What would be potential risks/difficulties e.g. for construction plans? Medical imaging (DICOM)?

.....

Logical Preservation

Questions / discussion (3):

- Which objects are most at risk? 
- Which objects are most difficult to preserve?
- How do we preserve entire business processes?
- If we loose significant properties with a strategy, what is the impact on authenticity? Can we use a “changed” object?
- What is the difference to systems engineering?

Why do we need Digital Preservation


3 levels of threat / preservation

1. Bit rot – physical preservation / bit preservation
Physical Layer: how to keep the 0's and 1's
2. Object formats – logical preservation
Logical Layer: how to remain able to open a file, run a program
3. Authenticity, interpretability – semantic preservation
Semantic Layer: how to ensure we can understand/interpret data correctly

- How? What can we do?

.....

Semantic preservation

- Threats at semantic level
 - meaning of terms change: city names, ... 
 - measurement scales, sensor sensitivity, ...change
 - interpretation of facts change: alcohol levels, ...
- Rather long-term, but subtle to notice
- Consider context of objects
 - purpose, setting, limitations, cultural context, related objects, ...

Semantic preservation

- Approaches / solutions:
 - Semantic enrichment
 - Metadata
 - Migration at semantic level
 - Documentation of context
 - Tracing of metadata
 - Document intended meaning / interpretation

Semantic preservation

Questions / discussion:

- How do we identify need for action?
- What is the risk of missing timely action?
- How do we solidly identify and document context?
- How can we implement semantic enrichment / semantic migration, ...?
- What about security issues?
- Is PDF save? PDF/A?
- Who is allowed to have access to which documents?
Who had access to them?
- Are differences in the communication protocol at an API level a problem of logical or semantic preservation?

.....

From Data to Processes



- Assume we know how to preserve data - **Is this sufficient?**
- Preserving data: Data Management Plans
 - describing data and context: provenance, authenticity, representation information,...
 - range of (ambiguous) definitions of context
 - But: mostly not actionable, not enforceable,...
 - BUT: data are (just) results of processes!
- Processes may be needed to
 - verify data
 - understand provenance
 - re-use process on new data
 - integrate data over time
- **Process curation instead of data curation!**

.....

What can go wrong?

- A lot.....
- Many times: trivial mistakes!
- But also more serious / conceptual issues
- From mistakes to actual fraud
- Overlap with security research, digital forensics, ...
- Roles and Responsibilities, Policies, ...

.....

What can go wrong?

- Ingest / Standardization:
 - Who is performing the initial migration?
 - Who is liable?
 - Who will need to manage any problems subsequently?

- Migration
 - Something added? E.g. Word -> TXT, Excel -> TXT
 - Something lost?

- What is a PDF file?
 - A malicious invoice...
 - A multi-purpose paper:
<https://bigdata.uni-saarland.de/publications/p1972-dittrich.html>

.....

What can go wrong?

- Dissemination

- (c.f. OAIS Model: AIP to DIP transformation)
- Decide which information to release in which format
- Consider metadata!

c.f. Supriya Adhatarao, Cédric Lauradoux: Exploitation and Sanitization of Hidden Data in PDF Files. arXiv:2103.02707, March 2021

Analyzing PDFs published by security agencies: metadata revealing weak links, less than 10% of agencies sanitized part of their documents, 65% of these still contained sensitive information)



What can go wrong?

- Collection Profiling: what is in the repository?



JHOVE2

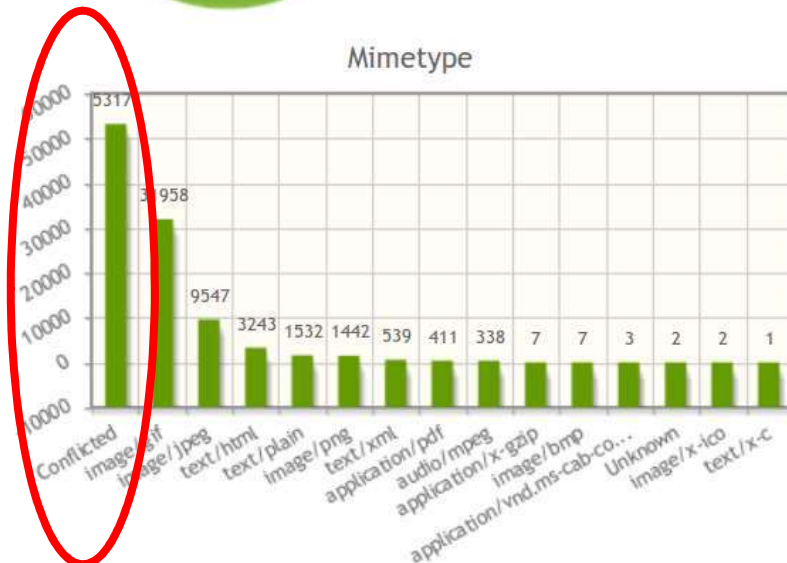


ffident

Droid



METADATA XTRACTOR



Overview

-
- What are the challenges in Digital Preservation? (recap)
 - What methods can we use to counter them?
 - What's the issue with file formats?
 - What do we have in our repository?
-

What Kind of File is This?

- By external characteristics (file extensions)
- By internal characteristics („magic number“, „signature“).

E.g. A TIFF file begins with ...

- Bytes 0-1
 - The byte order used within the file. Legal values are: “II” (4949.H) / “MM” (4D4D.H)
- Bytes 2-3
 - An arbitrary but carefully chosen number (**42**) that further identifies the file as a TIFF file.

What Kind of File is This?

- What's wrong with file extensions?
 - Not necessarily unique (e.g. wks)
 - Granularity not sufficient
 - Can be altered by users
 - Formats vs. Format profiles
 - PDF is not **one** format
 - DOC is not **one** format
 - TIFF is not **one** format
 - A lot of things can go wrong: by coincidence or maliciously!
 - Word -> TXT, Excel -> TXT
 - Standardization: “We only use PDFs, we've no problem”
 - Not all PDFs are created equal...
 - A PDF file?
- <https://bigdata.uni-saarland.de/publications/p1972-dittrich.html>

What's Wrong with MIME Types?

- Insufficient depth of detail
 - No requirements regarding syntax and semantic description
 - No requirement for complete disclosure, especially of proprietary formats
- Insufficient granularity
 - Both tiled RGB GeoTIFF with LZW and striped bi-tonal TIFF-FX with Group 4 are typed as “image/tiff”
 - All of PDF 1.0 – 1.4, PDF/X-1, X-2, X-3, and PDF/A are typed as “application/pdf”
 - These variants might require radically different workflows

Why Do We Need a Format Registry?

- Knowledge base of file format representation information
 - properties,
 - what do they mean?
 - how to read them?
 - supporting software
- Unification of vocabulary (entity names and mappings)
- A (single?) access point to various information about formats through a common API

File Format Registries

- PRONOM
 - <http://www.nationalarchives.gov.uk/pronom/>
- Global Digital Format Registry (defunct)
 - http://library.harvard.edu/preservation/digital-preservation_gdfr.html
- Unified Digital Format Registry (UDFR) (defunct)
 - <http://www.udfr.org/>
- Sustainability of Digital Formats Planning for Library of Congress Collections
 - <https://www.loc.gov/preservation/digital/formats/index.shtml>
- FileExt
 - <http://filext.com>



Details for: Microsoft Word for Windows Document 97-2003

Go to: [Summary](#) | [Documentation](#) > | [Signatures](#) > | [Compression](#) > | [Character encoding](#) > | [Rights](#) > | [Reference files](#) > | [Properties](#) >

Summary

Name	Microsoft Word for Windows Document
Version	97-2003
Other names	Microsoft Word for Windows Document (97-XP)
Identifiers	MIME: application/msword Apple Uniform Type Identifier: com.microsoft.word.doc PUID: fmt/40
Family	
Classification	Text (Wordprocessed)
Disclosure	None
Description	With the release of Word 97, Microsoft revised the native binary word processing format, which is based on its generic OLE2 Compound Document Format. The format is proprietary and Microsoft does not make details of its structure public. The information here is derived primarily from OpenOffice.org's reverse-engineered documentation of the format and should not therefore be regarded as definitive. A Word document is stored as a 'WordDocument' stream within a Compound Document Format file. The format remained unchanged with the releases of Word 2000, 2002 and 2003.
Orientation	Binary
Byte order	Little-endian (Intel)
Related file formats	Has priority over OLE2 Compound Document Format Is subsequent version of Microsoft Word for Windows Document (6.0/95) Is subtype of OLE2 Compound Document Format

Registry Content

- Descriptive information
- Identifiers
 - MIME
 - Pronom Unique Identifier (PUID)
- Relationships to formats
- Technical environment
- References and links
- Risk factors

.....

Registry Use Cases

- Identification
 - “I have a digital object; what format is it?”
- Validation
 - “I have an object purportedly of format F ; is it?”
- Transformation
 - “I have an object of format F , but need G ; how can I produce it?”
- Characterization
 - “I have an object of format F ; what are its significant properties?”
- Risk assessment
 - “I have an object of format F ; is it at risk of obsolescence?”
- Delivery
 - “I have an object of format F ; how can I render it?”

Identification Tools

- DROID (Digital Record Object Identification)
 - relies on PRONOM
 - The National Archives, UK
- JHOVE
 - JSTOR/Harvard Object Validation Environment
 - Validation and characterisation
- FITS (File Information Tool Set)
- Apache Tika: file content analysis
- veraPDF: PDF/A validation
<http://verapdf.org/home/>



ffident

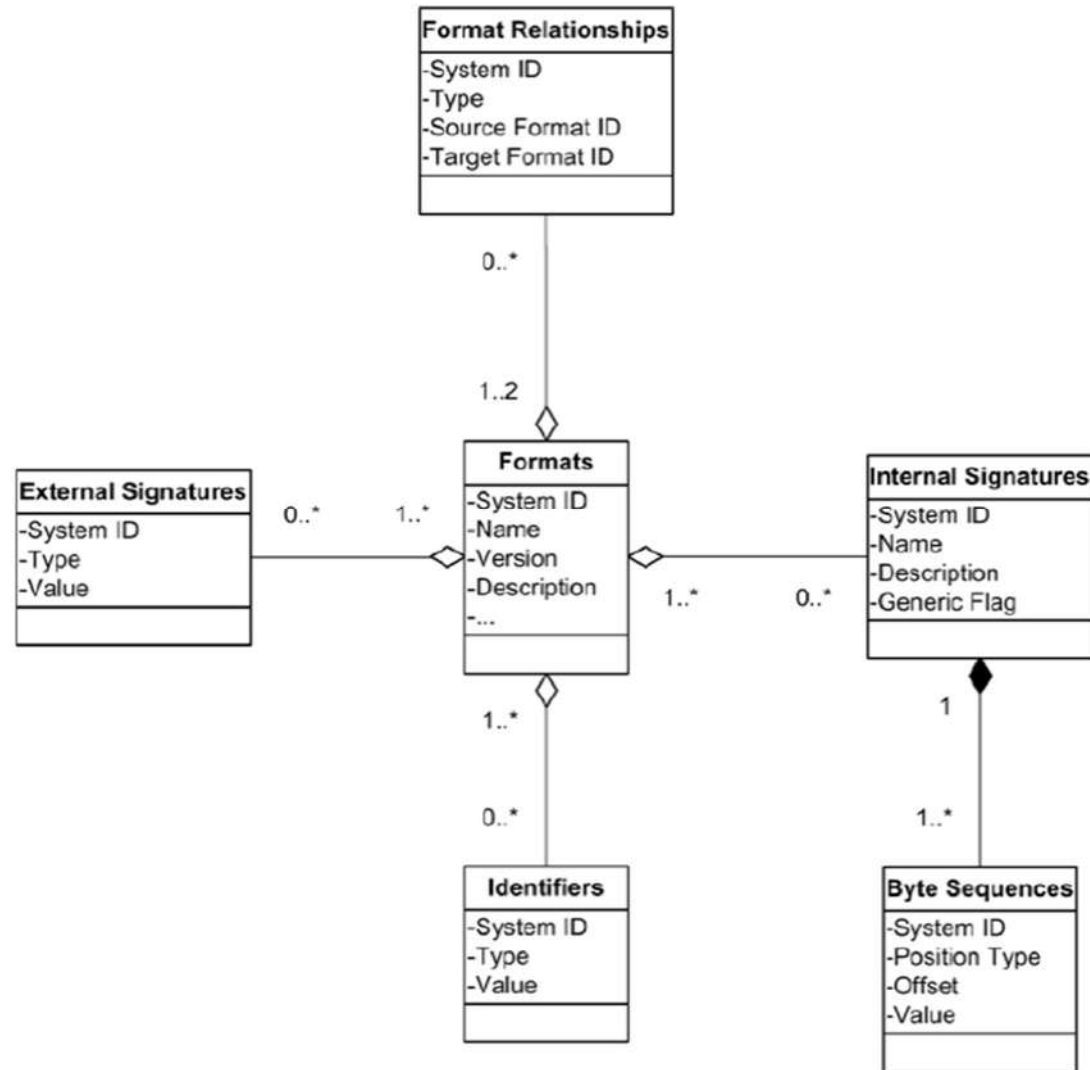


Droid



Signatures in DROID

- External signatures
 - File extensions
- Internal signatures
 - Format indicators in the bitstream
 - Byte sequences



.....



JHove

- JSTOR/Harvard Object Validation Environment
- Modular and extensible Java-based architecture
 - Image modules: GIF, JPEG, JPEG2000, TIFF
 - Document modules: ASCII,HTML,PDF, UTF-8, XML
 - ...
- Three functions
 - Identification
 - Validation
 - Characterisation
- JHove2
 - Identification and validation
 - Feature extraction
 - Policy based assessment
 - Able to handle complex objects

.....



The TIFF Module...

- Tagged Image File Format (TIFF) raster images TIFF 4.0, 5.0, and 6.0 [[TIFF 4.0](#), [TIFF 5.0](#), [TIFF 6.0](#)]
- Baseline 6.0 Class B, G, P, and R [[TIFF 6.0](#)]
- Extension Class Y [[TIFF 6.0](#)]
- TIFF/IT (ISO 12639:2003) [[TIFF/IT](#)] File types CT, LW, HC, MP, BP, BL, and FP, and conformance levels P1 and P2
- TIFF/EP (ISO 12234-2:2001) [[TIFF/EP](#)]
- Exif 2.0, 2.1 (JEIDA-49-1998), and 2.2 (JEITA CP-3451) [[Exif 2.1](#), [Exif 2.2](#)]
- GeoTIFF 1.0 [[GeoTIFF](#)]
- TIFF-FX (RFC 2301) [[TIFF-FX](#)]
- Profiles C, F, J, L, M, and S
- Class F (RFC 2306) [[Class F](#), [RFC 2306](#)]
- RFC 1314 [[RFC 1314](#)]
- DNG (Adobe Digital Negative) [[DNG](#)]



FITS - File Information Tool Set

✦ Main features:

- Consolidates output
- Can include raw output
- Configurable/Extendable

✦ FITS includes:

- Droid
- Metadata Extra
- Jhove
- Exiftool
- FFident
- File Utility

.....



Conflicts

3 types of conflicts:

1. Inconsistent property naming, e.g: *image_width* and *imagewidth*
2. Competing characterisation results, e.g: tool1 identifies a file as *plain text*, but tool2 identifies the file as *PDF*
3. Close, but not the same property values, e.g: *application/xhtml+xml* vs. *application/xml*.

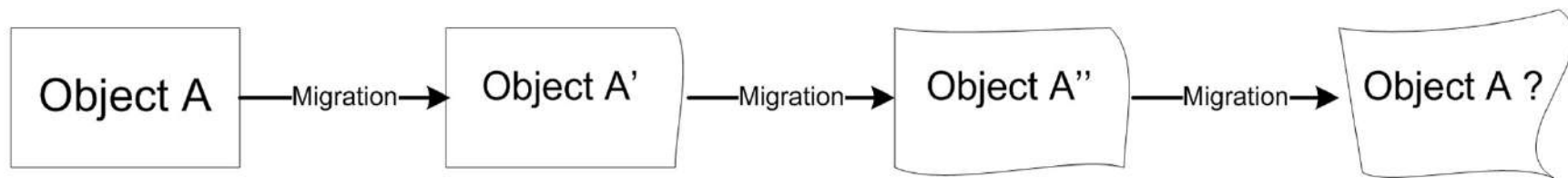
Validation

- A digital object is **well-formed** if it meets the purely syntactic requirements for its format.
- An object is **valid** if it is well-formed and it meets additional semantic-level requirements.

- Validation use cases:
 - "I have an object that purports to be of format F ; is it?"
 - "I have an object of format F ; does it meet profile P of F ?"
 - "I have an object of format F and external metadata about F in schema S ; are they consistent?"

Validating a migrated image

- Yes, it's in JPEG 2000 format
- Yes, it's well-formed
- Yes, it's valid
- Yes, it still has the same dimensions
- But is it still the same image?
- We need more characterisation.

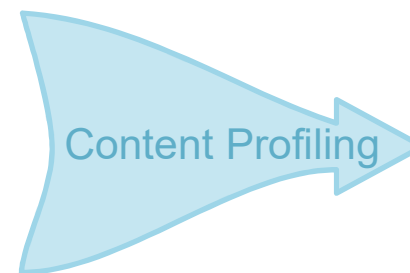
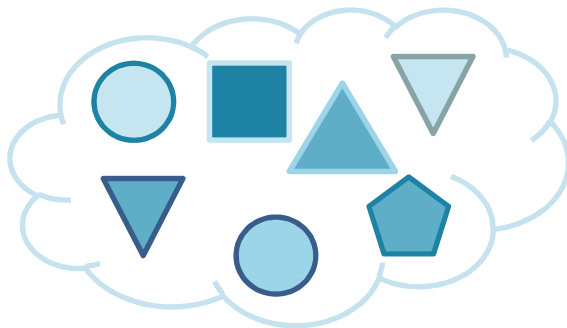


.....

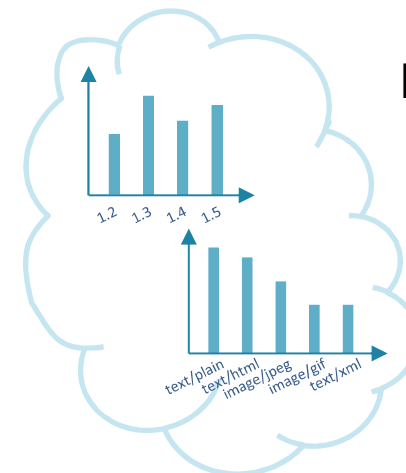
What is Content Profiling

- Better understand your content
- Reveal risks and opportunities
- Part of Preservation Planning
 - Planning and Watch, <http://bit.ly/scape-suite>

Your content



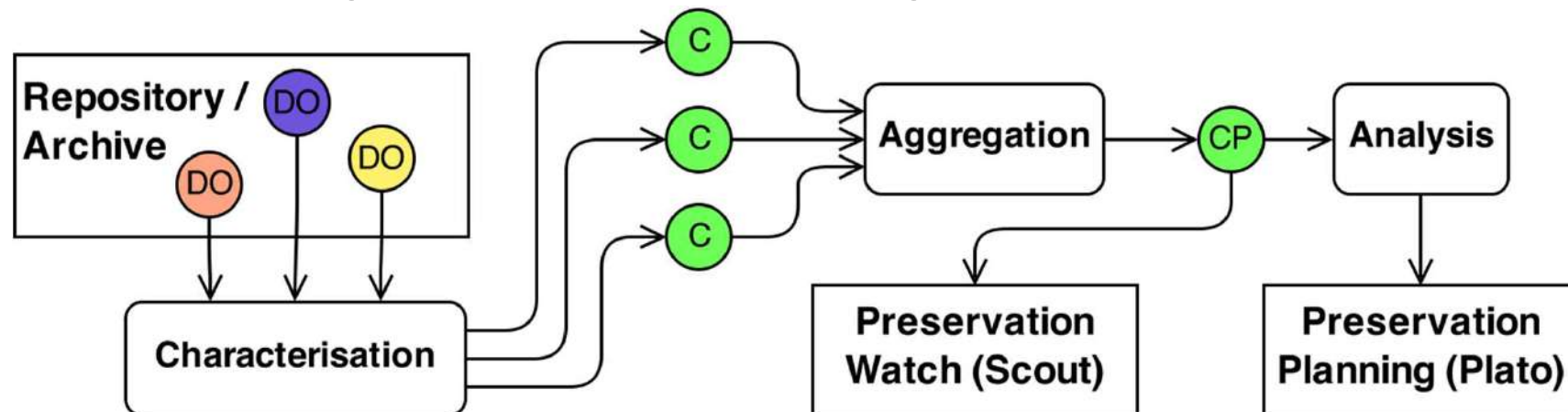
Report



.....

Content Profiling in Details

- A way of getting control over data
 - Decision support
- Consists of:
 - Characterization
 - Aggregation
 - Analysis
 - Reporting / Use for decision making



Aggregation

- Provides an overview of the content
 - Distributions of characteristics
 - Statistics (size, min, max, avg...)
- Data sizes grow dramatically
- Heterogeneity of data
- No universal characterisation tool
 - Combination of such tools

Analysis

- In-depth research into your content
 - Drilling down
 - Filtering
 - SQL-like queries
- Representative samples generation
 - Based on metadata
 - Outlier detection
 - Stratification across
 - File type,
 - Size,
 - Time, or
 - Any other relevant property

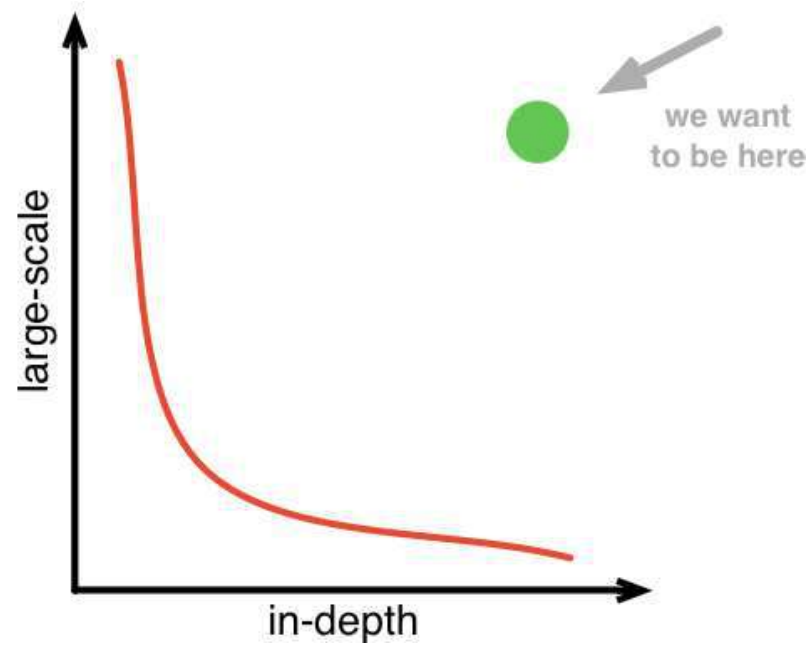
.....

Challenges

- Lack of:
 - Trustworthy tools for characterisation
- Depth
 - Address heterogeneity of data
 - Rise awareness of content properties
 - Combine several characterisation tools
- Quality
 - Conflicts due to combination of characterisation results
 - Resolve conflicting metadata

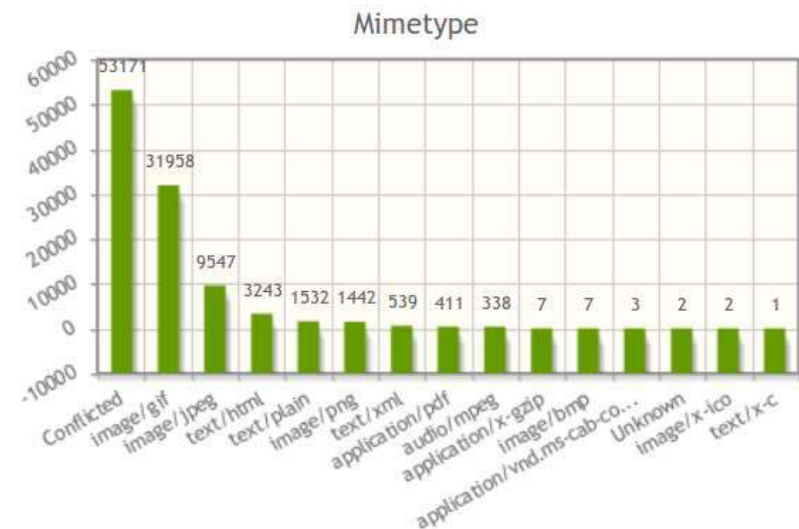
Challenges

- Scale
 - Effectively analyze substantial amount of metadata
 - Large-scale approaches for content profiling



- C3PO – Clever, Crafty, Content Profiling of Objects
- Reads and analyzes information from FITS
- Support large-scale database solutions for aggregation and analysis of characterisation metadata
 - MongoDB, <http://www.mongodb.org>
 - HBase, <http://hbase.apache.org>
- Aggregation-only mode
 - Useful to fast and explorative generation of a content profile
 - Statistics calculation using predefined filters
 - Single read of data, without computationally expensive ingest and further analysis

- Uses characterisation results
- Interface to support other characterisation tools
- Deeper content analysis with interactive visuals through a web-app
- Representative sampling
- Open-source
 - <http://ifs.tuwien.ac.at/imp/c3po>



- Stored metadata property mapping to the existing vocabulary, Planning and Watch Ontology
 - <http://purl.org/DP/quality/asures>

 compression algorithm (Individual) Definition

Definition
The **URI** of this individual is `http://purl.org/DP/quality/asures#118`

compression algorithm	<code>http://www.w3.org/2004/02/skos/core#prefLabel</code>	compression algorithm
compression algorithm	<code>http://purl.org/DP/quality#scale</code>	<code>http://purl.org/DP/quality/scales#FREETEXT</code>
compression algorithm	<code>http://purl.org/DP/quality#attribute</code>	<code>http://purl.org/DP/quality/attributes#39</code>

- Rule-based engine to resolve conflicts in characterization metadata
 - Drools, <http://www.jboss.org/drools>
- Preservation-specific rules

Target tool	Execution condition	Action
Droid, Exiftool, all	Droid and Exiftool identify a file as “Microsoft Powerpoint Presentation”	Ignore format identifications by other tools
Jhove, all	Jhove reports “text/html” mimetype, other tools report “application/xhtml+xml”	Ignore the “text/html” mimetype provided by Jhove

Overview

-
- What are the challenges in Digital Preservation? (recap)
 - What methods can we use to counter them?
 - What's the issue with file formats?
 - What do we have in our repository?
-

Digital Preservation - Summary

- Is a complex task
- Requires a concise understanding of the objects, their intellectual characteristics, the way they were created and used and how they will most likely be used in the future
- Requires a continuous commitment to preserve objects to avoid the „digital dark hole“
- Requires a solid, trusted infrastructure and workflows to ensure digital objects are not lost
- Is essential to maintain electronic publications & data accessible
- Will become more complex as digital objects become more complex
- Needs to be defined in a preservation plan

Questions / Discussion:

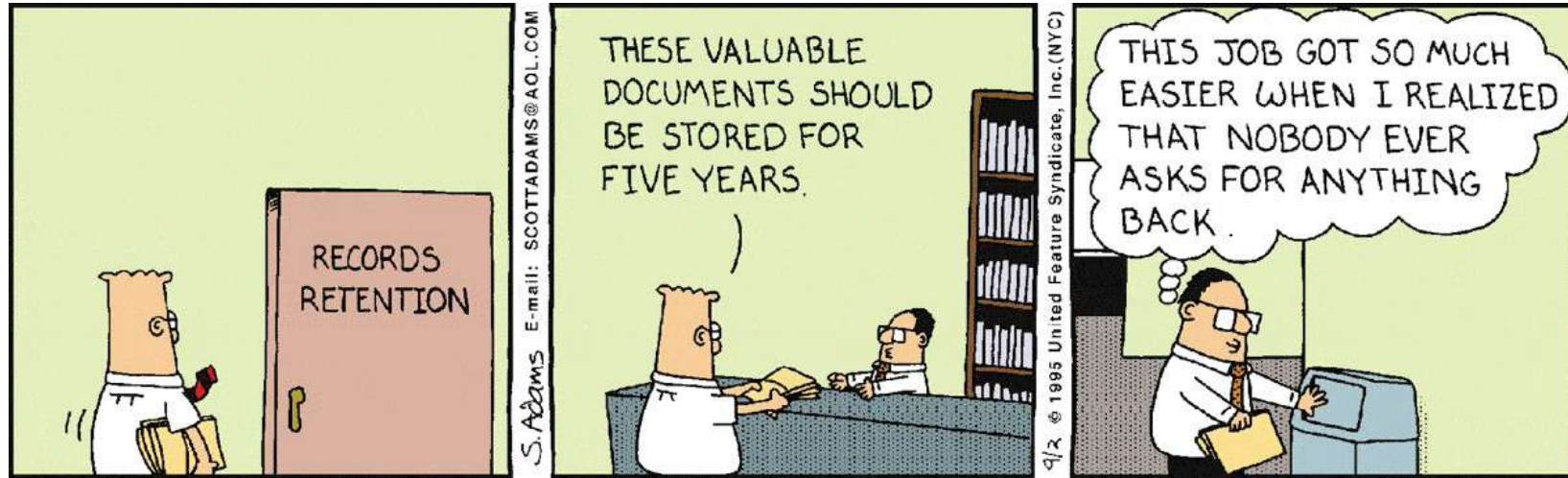
- At what levels are digital objects threatened?
- What are the time intervals at each level?
- How can we identify objects at risk?
- What can we do to mitigate the risk?
- How can we recover if mitigation fails / is missed?
- What competences do we need?
- How would a training/education program look like?
- How do we know if somebody is doing a good job at DP?

Current Issues

- Atomic file formats, stability of file formats
 - What are the atomic building blocks of information?
 - Can we split information objects?
 - Can we synthesize them? - Help for benchmarking?
- Scalability, Semantics
- Digital forgetting
 - how to decide what to keep and what to forget?
 - keep all? just storage? how to find? utilize? understand?
- Sustainable Systems Engineering
 - How can we build preservation-ready systems?
 - How to integrate DP-considerations into software engineering?
- Costs: what does DP cost?
 - cost factors?
 - How to model? evaluate?

.....

Thank you!

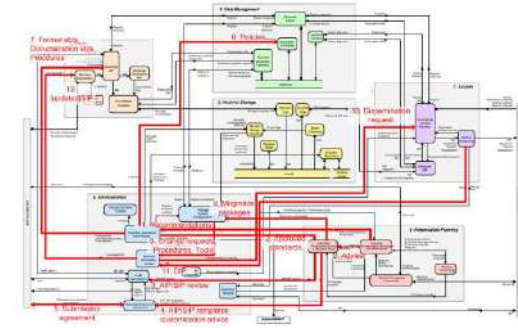
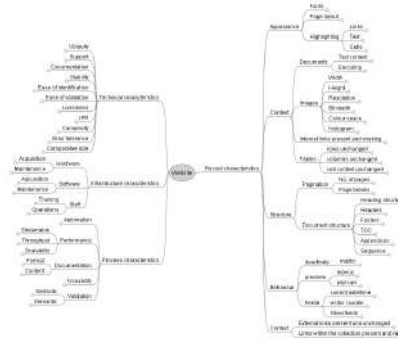
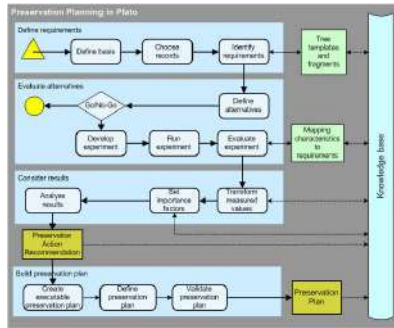


Source: <http://dilbert.com/strip/1995-09-02>

<http://www.ifs.tuwien.ac.at/dp>

.....

Thank you!



<http://www.ifs.tuwien.ac.at/dp>

ID	Name	Value	Unit	Priority	Weight	Link
1	1.1.1.1	1000	MB	High	1.0	
2	1.1.1.2	2000	MB	Medium	0.5	
3	1.1.1.3	3000	MB	Low	0.2	
4	1.1.1.4	4000	MB	High	1.0	
5	1.1.1.5	5000	MB	Medium	0.5	
6	1.1.1.6	6000	MB	Low	0.2	
7	1.1.1.7	7000	MB	High	1.0	
8	1.1.1.8	8000	MB	Medium	0.5	
9	1.1.1.9	9000	MB	Low	0.2	
10	1.1.1.10	10000	MB	High	1.0	



Digital Preservation

OAIS Reference Model

Andreas Rauber
Institute of Software Technology and Interactive Systems
Vienna University of Technology
<http://www.ifs.tuwien.ac.at/dp>

Outline

-
- Principles of the OAIS Model
 - Technical Overview
 - Functional Overview
 - Information Modell
 - Summary
-


OAIS and the role of NASA

- National Space Science Data Center
 - NASA's first digital archive
 - has gone through many technology changes since 1966
- Consultative Committee for Space Data Systems
 - International group of Space Agencies
 - developed a set of standards across disciplines
 - evolved into working group ISO TC 20/ SC 13 around 1990
 - TC20: Aircraft and Space Vehicles
 - SC13: Space Data and Information Transfer Systems

What's a reference model

- A Framework
 - to understand the relationship between significant entities in an environment
 - for the development of consistent standards or specifications to support this environment.
- A reference model
 - is based on a small number of unified concepts
 - is an abstraction of the core concepts, their relationships and interfaces within as well as external to the framework
 - can be used as a basis for training and to explain standards.

OAIS

- OAIS is a reference model 
- No design specification, no data model, no set of functional requirements!
- Describes elements and concepts that are relevant for a project
- Goal: determine, which parts of the reference model map to which subsystems, functions and responsibilities in a desired solution.

OAIS Sources of Information

- Reference Model for an Open Archival Information System (OAIS), ISO 14721:2012
- Blue Book, CCSDS 650.0-B-1, January 2002
- Pink Book, CCSDS 650.0-P-1.1, August 2009
- <http://public.ccsds.org/sites/cwe/rids/Lists/CCSDS%206500P11/Attachments/650x0p11.pdf>
- Slides based on Blue Book, Pink Book and:
 - Don Sawyer, Lou Reich: ISO Reference Model for an Open Archival Information System (OAIS) Tutorial Presentation, LOC, June 13 2003
- <http://ssdoo.gsfc.nasa.gov/nost/isoas/overview.html>

Outline

-
- Principles of the OAIS Model
 - Technical Overview
 - Functional Overview
 - Information Modell
 - Summary
-

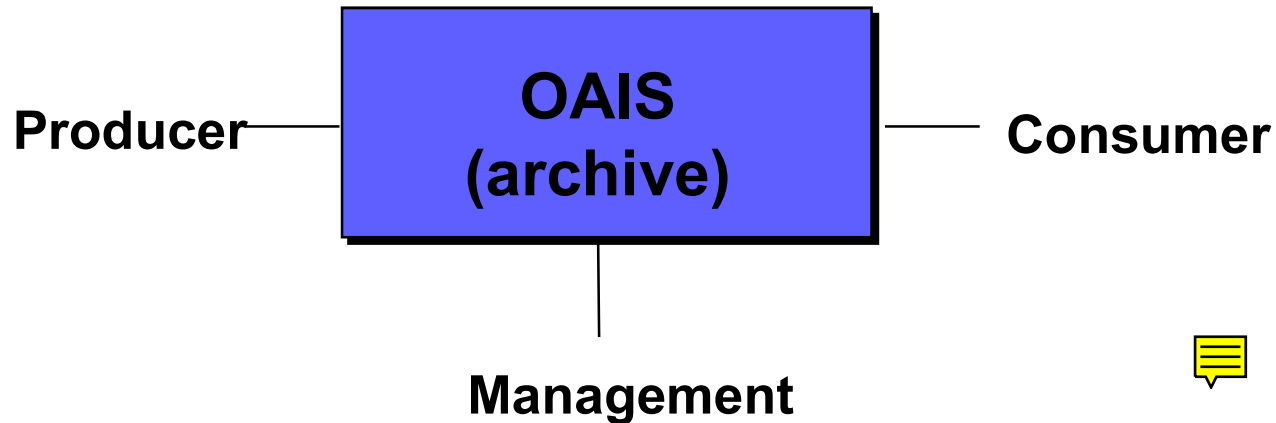
- Open
 - Reference Model standard(s) are developed using a public process and are freely available
- Information
 - Any type of knowledge that can be exchanged
 - Independent of the forms (i.e., physical or digital) used to represent the information
 - Data are the representation forms of information
- Archival Information System
 - Hardware, software, and people who are responsible for the acquisition, preservation and dissemination of the information

Purpose, Scope, and Applicability

- Framework for understanding and applying concepts needed for long-term digital information preservation
 - Long-term is long enough to be concerned about changing technologies
 - Starting point for model addressing non-digital information
- Provides set of minimal responsibilities to distinguish an OAIIS from other uses of 'archive'
- Framework for comparing architectures and operations of existing and future archives
- Basis for development of additional related standards
- Addresses a full range of archival functions
- Applicable to all long-term archives and those organizations and individuals dealing with information that may need long-term preservation
- Does NOT specify an implementation

.....

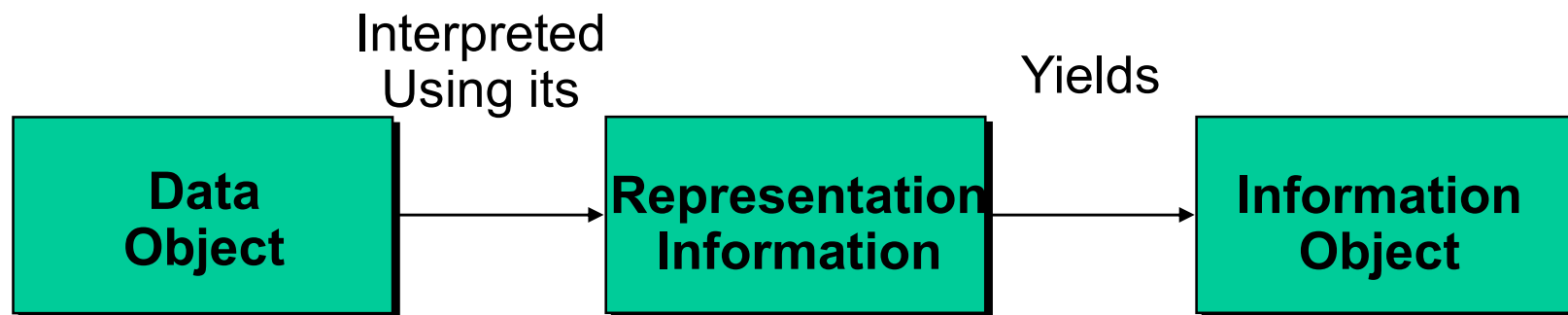
Model View of an OAIIS Environment



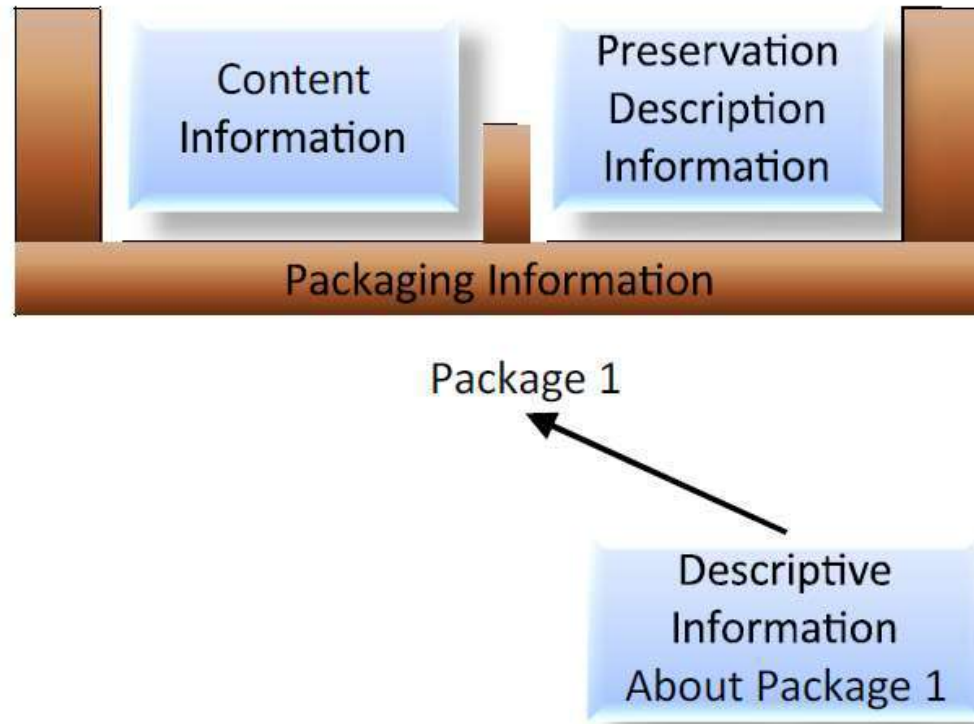
- Producer is the role played by those persons, or client systems, who provide the information to be preserved
- Management is the role played by those who set overall OAIIS policy as one component in a broader policy domain
- Consumer is the role played by those persons, or client systems, who interact with OAIIS services to find and acquire preserved information of interest

OAIS Information Definition

- Information is always expressed (i.e., represented) by some type of data
- Data interpreted using its Representation Information yields Information
- Information Object preservation requires clear identification and understanding of the Data Object and its associated Representation Information







Information Package Definition

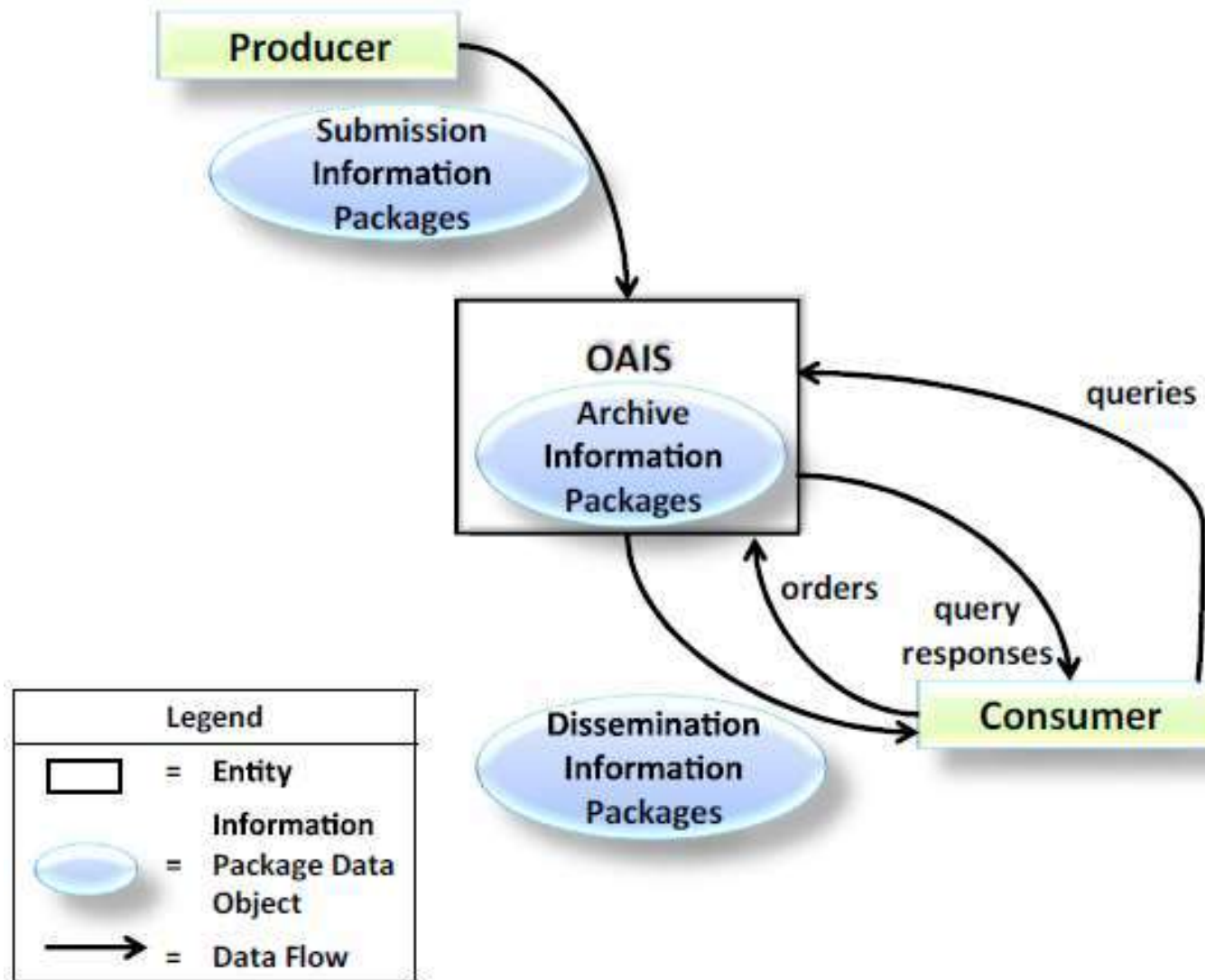


- An Information Package is a conceptual container holding two types of information
 - Content Information
 - Preservation Description Information (PDI)
 - Plus descriptive information

Information Package Variants

- **SIP: Submission Information Package** 
 - Negotiated between Producer and OAIS
 - Sent to OAIS by a Producer 
- **AIP: Archival Information Package** 
 - Information Package used for preservation
 - Includes complete set of Preservation Description Information (PDI) for the Content Information
- **DIP: Dissemination Information Package**
 - Includes part or all of one or more Archival Information Packages
 - Sent to a Consumer by the OAIS 

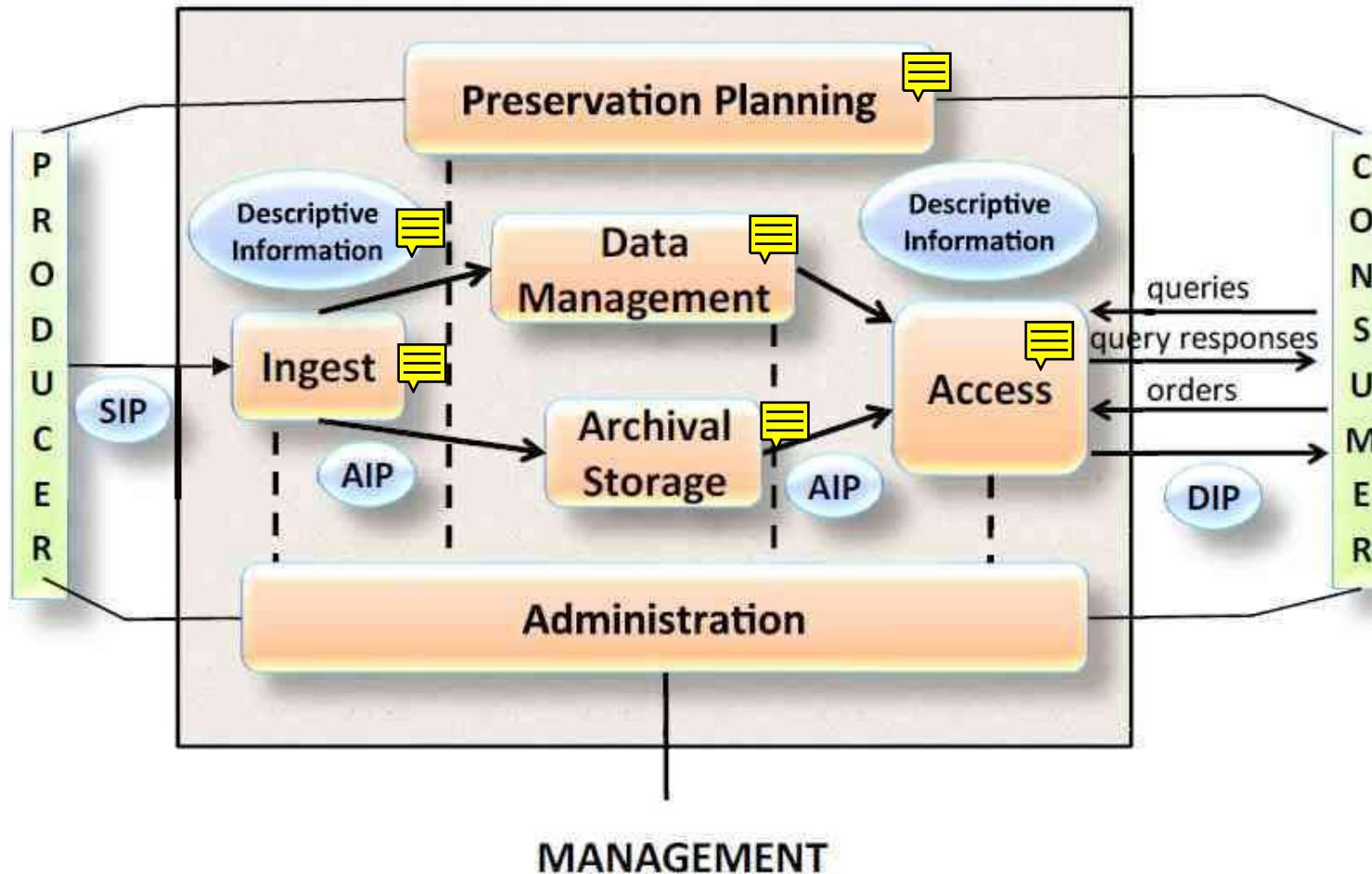
External Data Flow View



Outline

-
- Principles of the OAIS Model
 - Technical Overview
 - **Functional Overview**
 - **Information Modell**
 - Summary
-

Open Archival Information System: Six Functional Entities



SIP = Submission Information Package

AIP = Archival Information Package

DIP = Dissemination Information Package

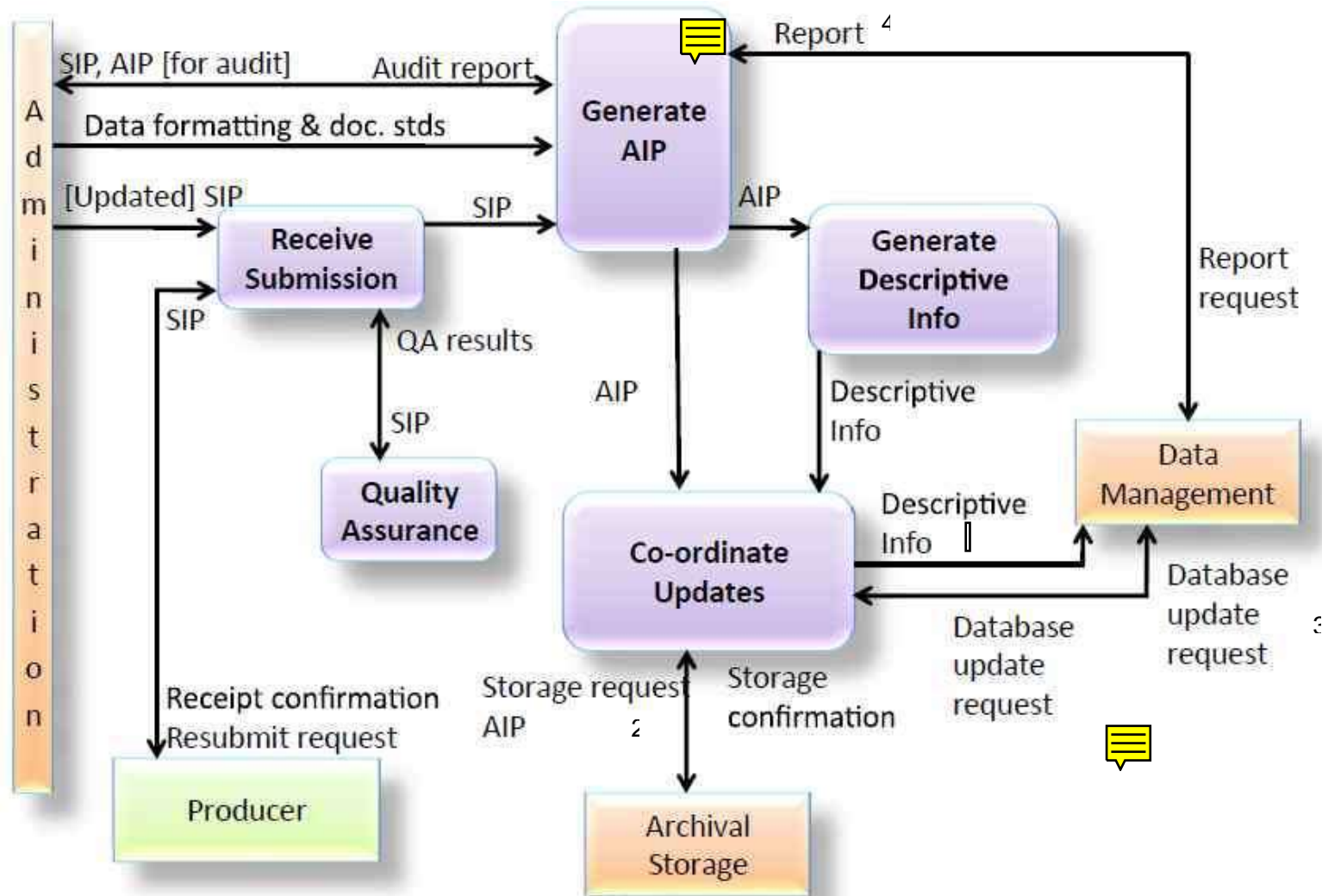
Functional Entities in an OAIS (1/2)

- **Ingest:** This entity provides the services and functions to accept Submission Information Packages (SIPs) from Producers and prepare the contents for storage and management within the archive
- **Archival Storage:** This entity provides the services and functions for the storage, maintenance and retrieval of Archival Information Packages
- **Data Management:** This entity provides the services and functions for populating, maintaining, and accessing both descriptive information which identifies and documents archive holdings and internal archive administrative data.

Functional Entities in an OAIS (2/2)

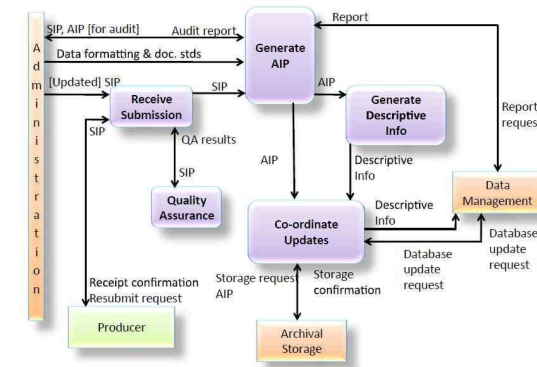
- **Administration:** This entity manages the overall operation of the archive system
- **Preservation Planning:** This entity monitors the environment of the OAIS and provides recommendations to ensure that the information stored in the OAIS remain accessible to the Designated User Community over the long term even if the original computing environment becomes obsolete.
- **Access:** This entity supports consumers in determining the existence, description, location and availability of information stored in the OAIS and allowing consumers to request and receive information products

Ingest Data Functions



Ingest Data Functions

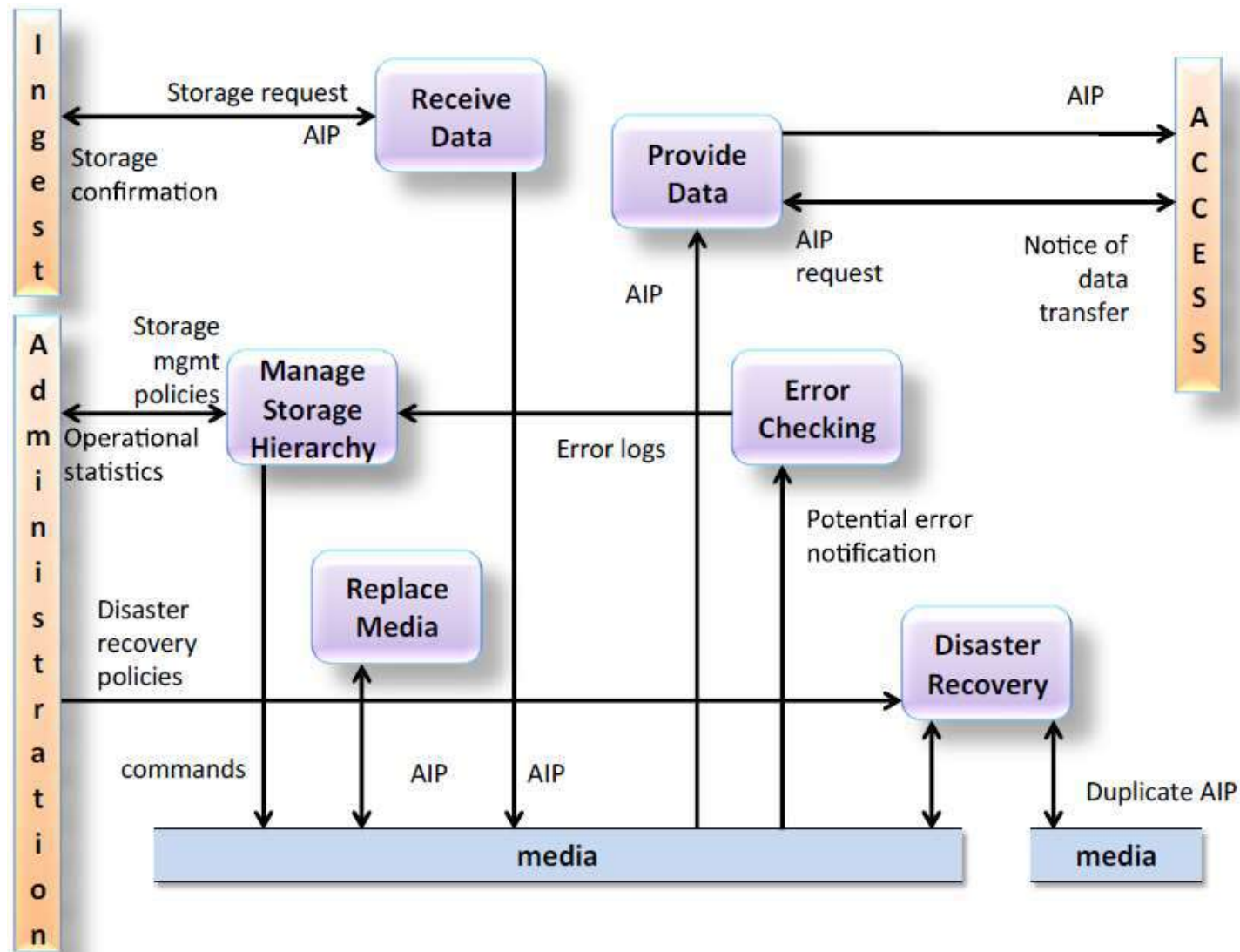
- Receive Submissions:
 - Staging Area for submissions
 - Confirmation on acceptance into staging area
- Quality Assurance
 - Validation of submission (CRC, logs, identity checks, media)
- Generate AIP
 - Transformation of SIPs into AIPs according to standard/policies (Transformation, Migration, Transcoding)
 - Forwarding of AIPs to Audit (Administration)
- Generate Descriptive Information
 - Collection / extraction of descriptive information on AIP for Data Management and Access Aids
- Coordinate Updates
 - Transfer of AIPs to Archival Storage
 - Confirmation -> Descriptive Information -> Data Management



Ingest Data Functions

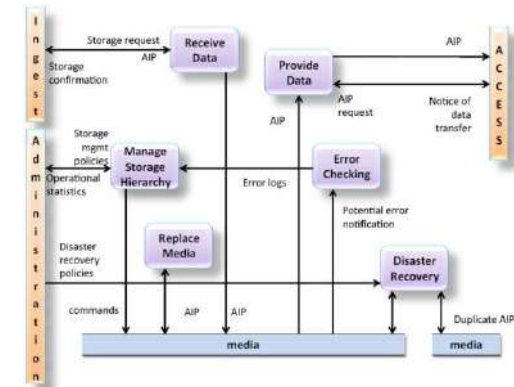
- Remember: what can go wrong?!
- Ingest / Standardization:
 - Who is performing the initial migration?
 - Who is liable?
 - Who will need to manage any problems subsequently?
- Migration
 - Something added? E.g. Word -> TXT, Excel -> TXT
 - Something lost?
- What is a PDF file?
 - A malicious invoice...
 - A multi-purpose paper:
<https://bigdata.uni-saarland.de/publications/p1972-dittrich.html>

Archival Storage Functions

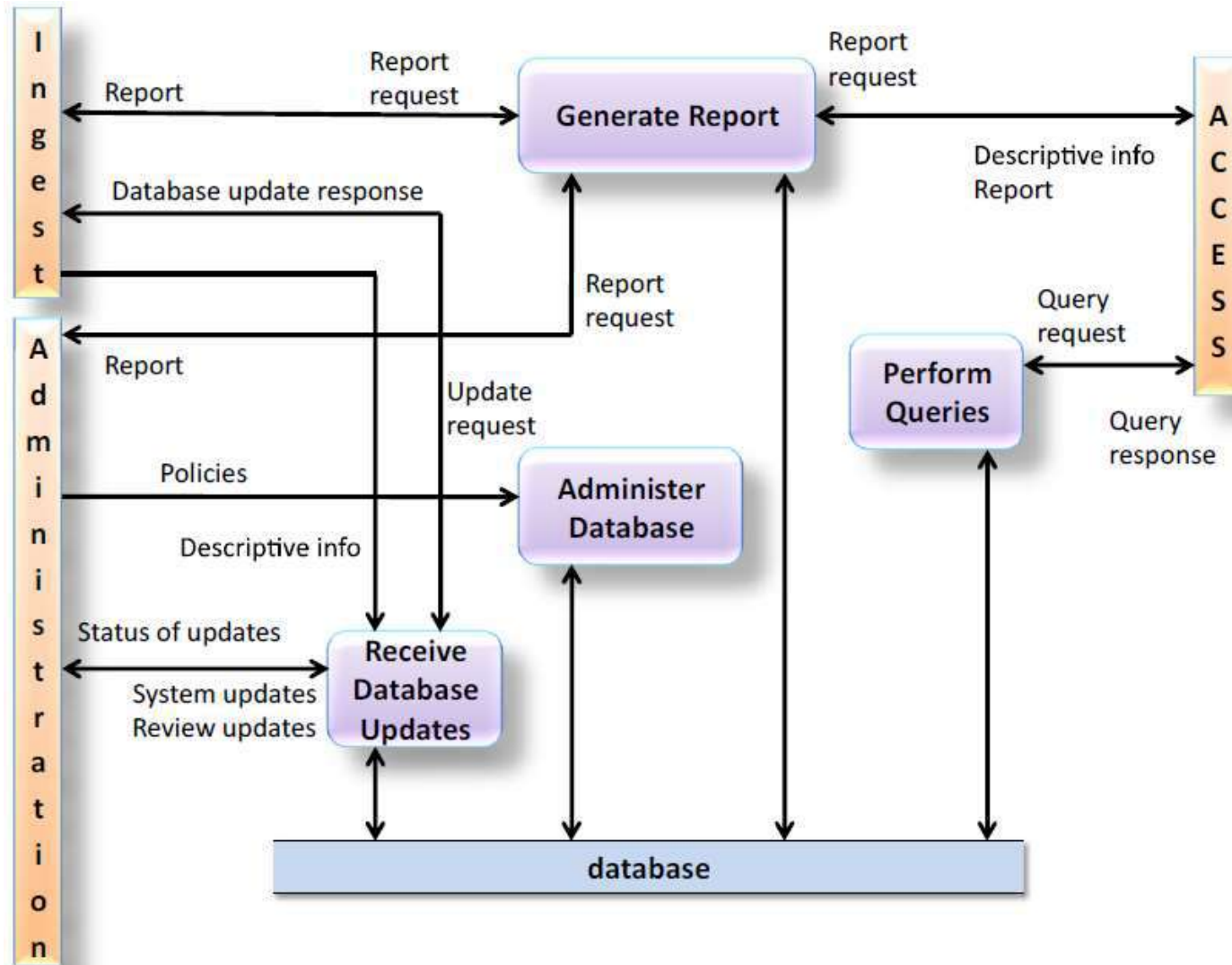


Archival Storage Functions

- Receive Data:
 - accept storage request for AIP
 - Decides on storage location, media
 - Returns confirmation messages
- Manage Storage Hierarchy
 - Management of storage according to policies
 - Monitoring of error messages, operational statistics
- Replace Media
 - Reproduction of AIPs over time (no changes of content or Preservation Description Information, only Packaging Information – other changes need to go via Administration)
- Error Checking
 - PDI Fixity Information (CRCs, error-correcting codes, ...)
- Disaster Recovery
 - Duplicating of storage media content (back-up)
 - Transport to physically separated location
- Provide Data
 - Generate copies of AIPs for Access



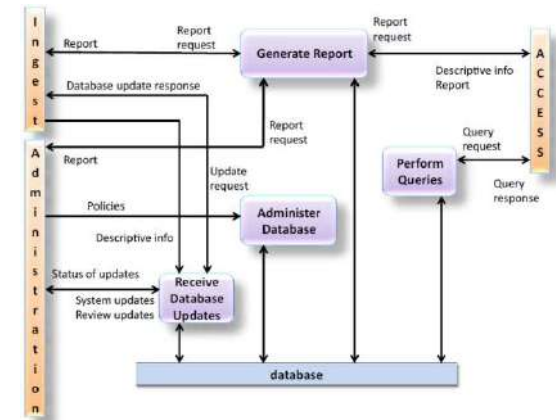
Data Management Functions



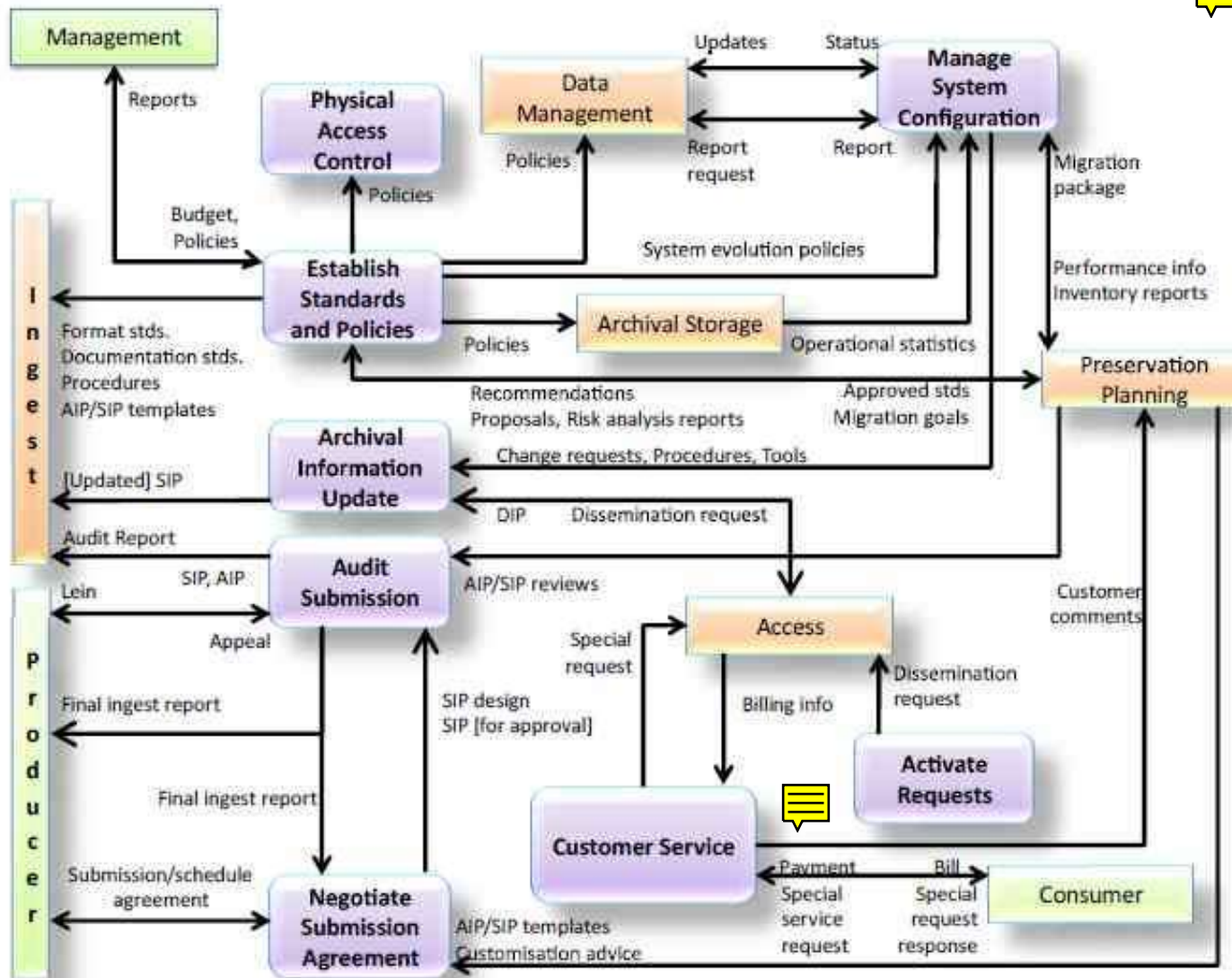
.....

Data Management Functions

- Administer Database
 - Integrity of DB for Descriptive Information and system information
- Perform Queries
 - Processing of queries from Access
- Generate Reports
 - Reports for Ingest, Access, Administration
- Receive DB Updates
 - Add/delete/modify information in Management DB
 - Ingest: new AIPs, Administration: updates

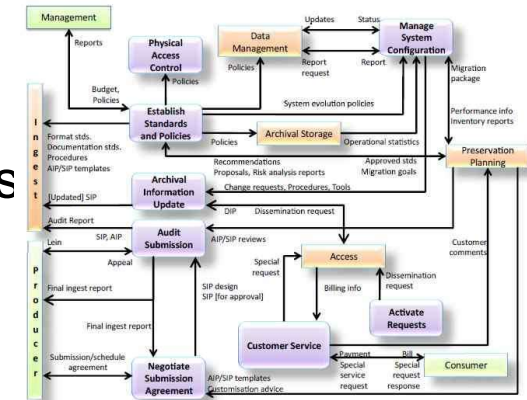


Administration Functions



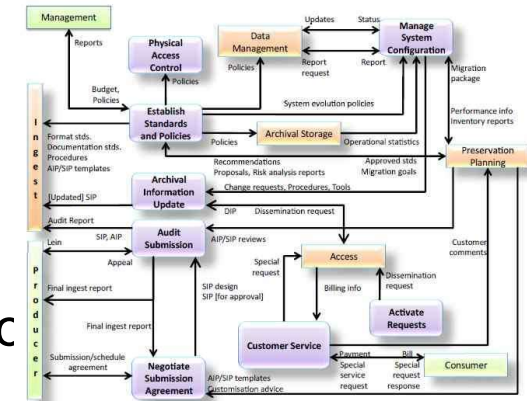
Administration Functions

- Negotiate Submission Agreement
 - Contracts with Producers, submission procedures
- Manage System Configuration
 - System evolution, monitoring
 - Provide information for Policies
- Archival Information Update
 - Update content of the archive: Modifying DIPs and Re-Submission -> Migration
- Establish Standards and Policies
 - Budget, Standards, Policies
- Audit Submission
 - Analyse whether SIPs and AIPs conform to policies and regulations
 - Verifia Representation and Package Information

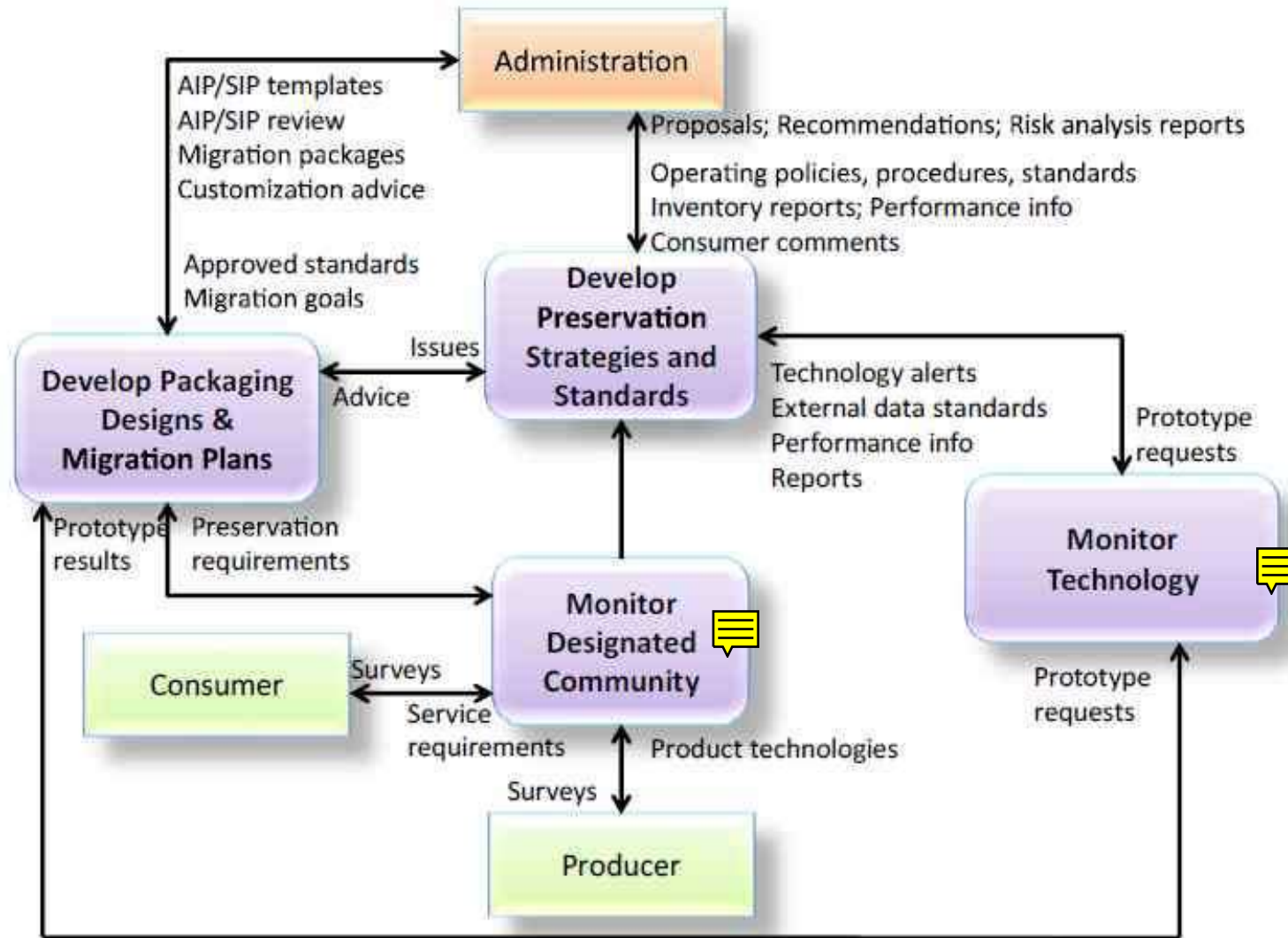


Administration Functions

- **Activate Requests**
 - Protocol of triggered queries/evaluations
 - Periodic checks/queries to archive to verify cc
 - Ordering data/reports periodically
- **Customer Service**
 - Manage customer accounts
 - Collect costs from Access, create invoices for customers

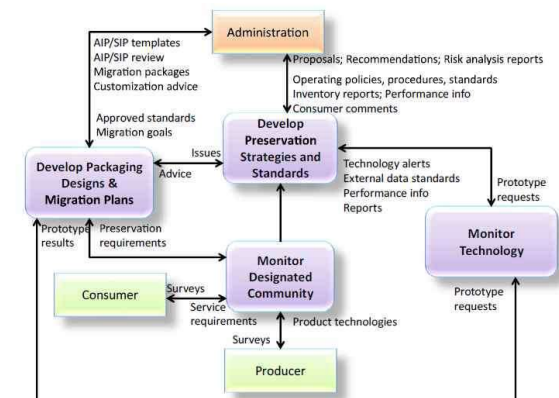


Preservation Planning Functions

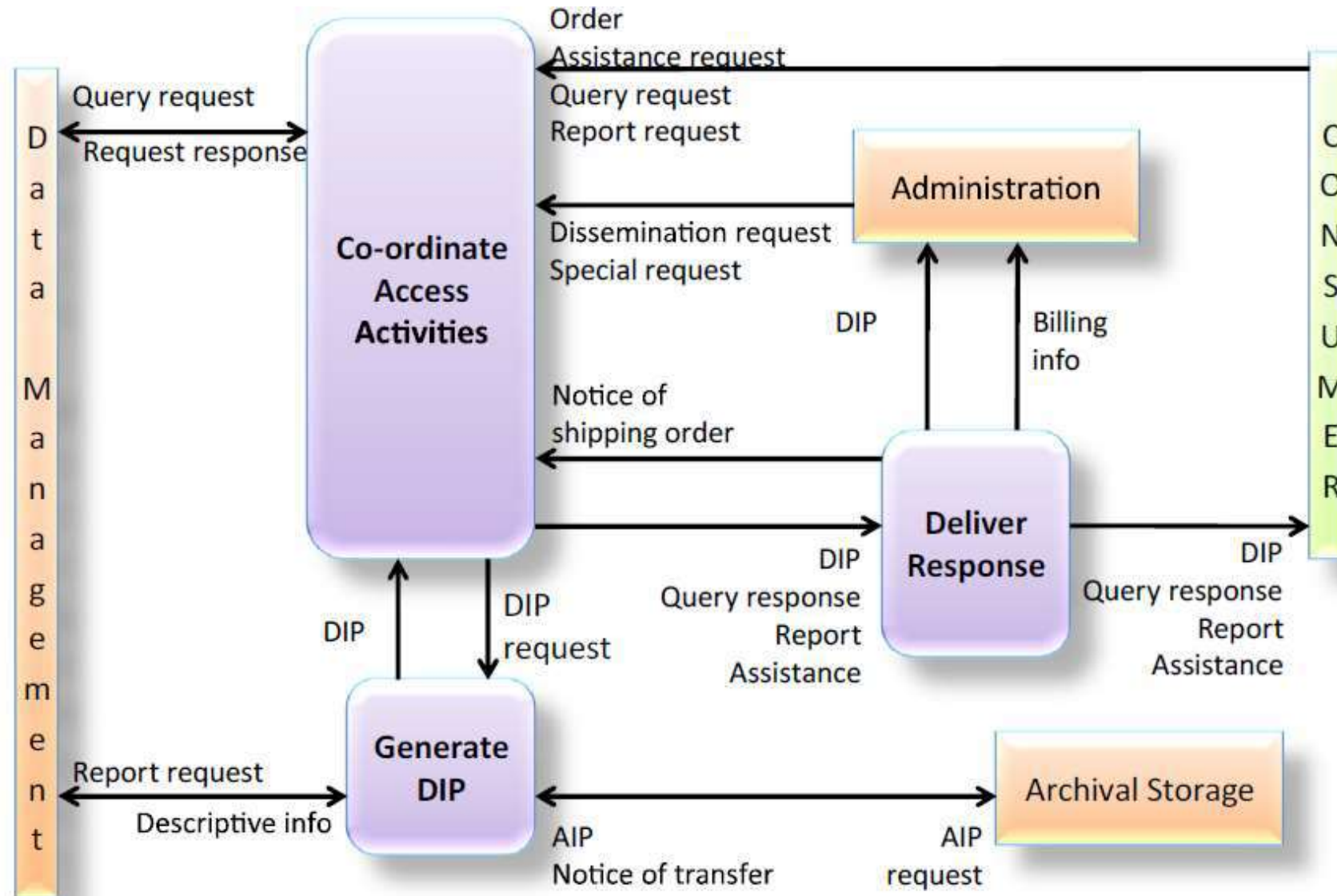


Preservation Planning Functions

- Monitor Designated Community
 - Interaction with Producer and Consumer
- Monitor Technology
 - Technology evolution: HW, SW, Formats
- Develop Preservation Strategies and Standards
 - Strategies, trend analysis
- Develop Packaging Designs and Migration Plans
 - Migration paths, tools
 - Create Preservation Description Information



Access Functions

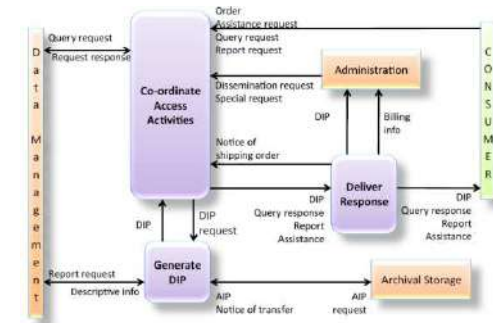


Access Functions

- Coordinate Access Activities
 - User interface, authorization
 - 3 types of Requests:
 - Queries to Data Management for Result Set
 - Order for Data Management and Archival Storage
 - Dissemination Requests by Administration for Archival Information Update

- Generate DIP
 - Get data from Archival Storage into Staging Area
 - Get Descriptive Information from Data Management
 - Apply processes to transform AIPs into a suitable DIP depending on query / consumer

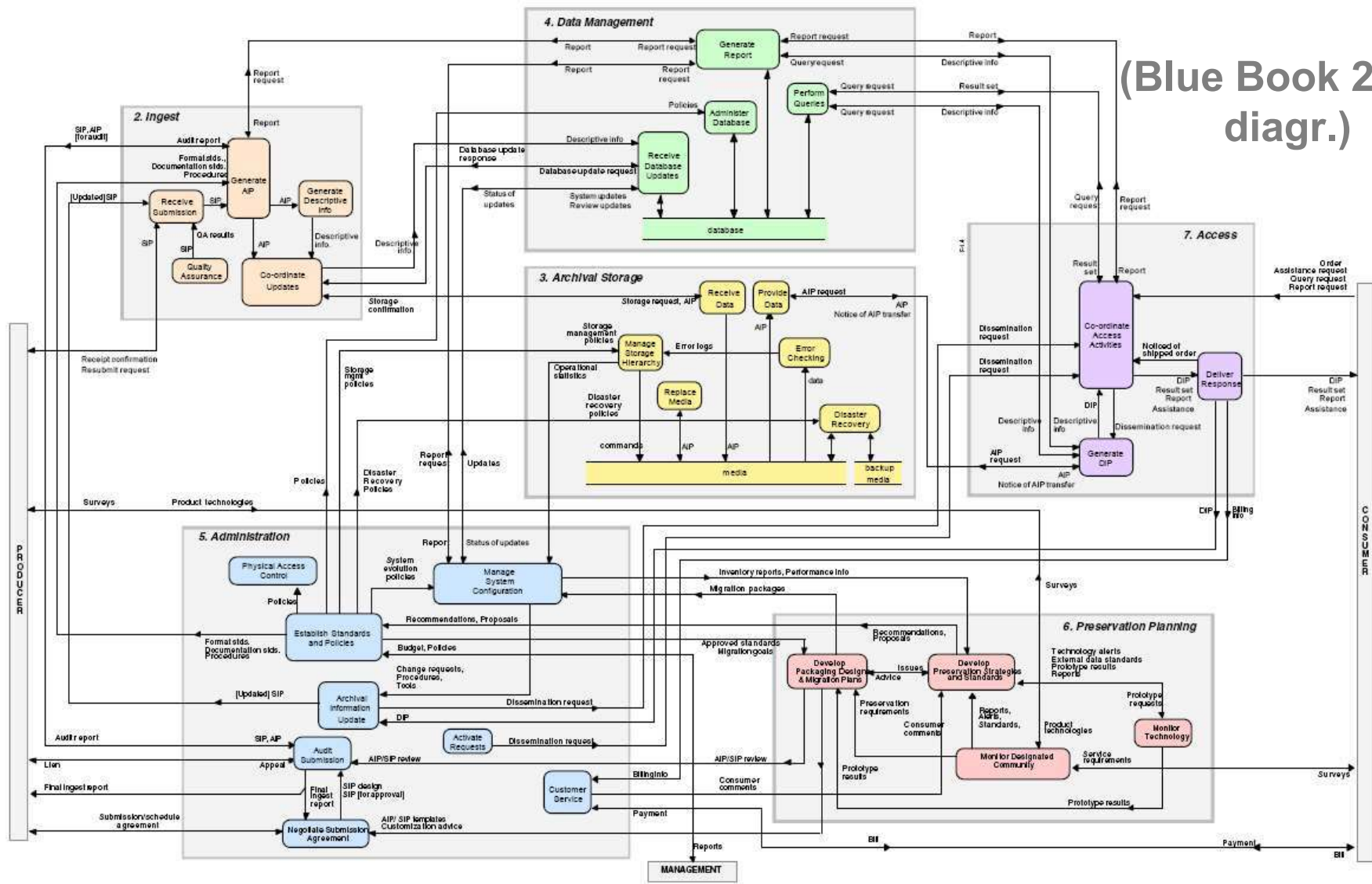
- Deliver Response
 - On-line and off-line responses
 - Forward results



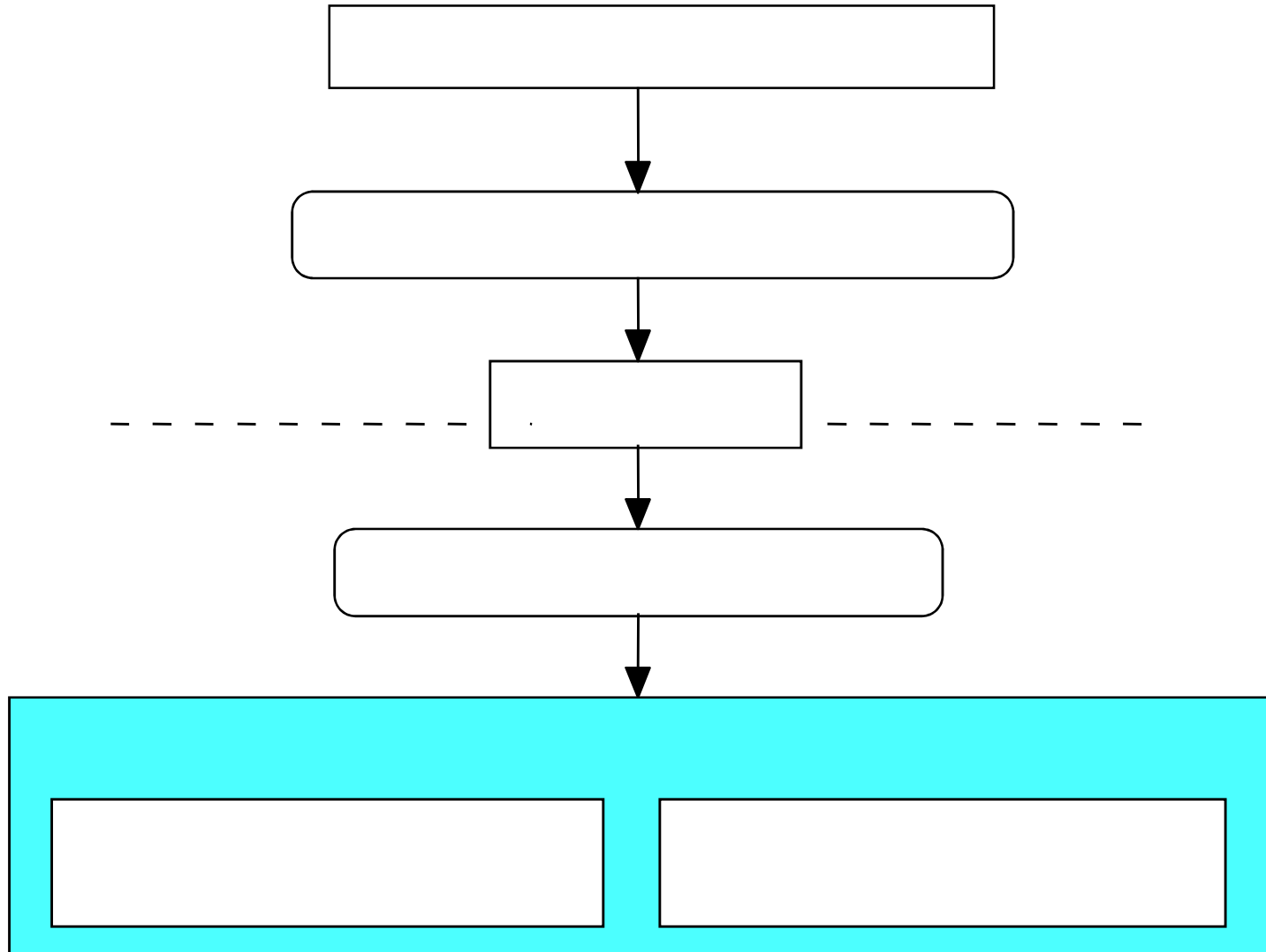
Access Functions

- AIP – DIP conversion
 - Decide which information to release in which format
 - Remember: Consider metadata!
c.f. Supriya Adhatarao, Cédric Lauradoux: Exploitation and Sanitization of Hidden Data in PDF Files. arXiv:2103.02707, March 2021
(Analyzing PDFs published by security agencies: metadata revealing weak links, less than 10% of agencies sanitized part of their documents, 65% of these still contained sensitive information)

OAIS Composite Functional Entities



Migration Context



Digital Migration

Digital Migration is defined to be the transfer of digital information, while intending to preserve it, within the OAIS.

- Focus on preservation of the full information content
- New information implementation replaces the old
- OAIS has full control and responsibility over all aspects of the transfer

Migration Motivators


- Motivators driving digital migrations
 - Media Decay
 - Often this is superseded by escalating media drive maintenance costs
 - Increased Cost Effectiveness
 - More cost-effective media types with higher volumes and lower drive maintenance costs
 - New User/Consumer Service Requirements
 - New formats more compatible with user's technology and applications
 - Proprietary software evolution
 - New software versions used to 'upgrade' formats of the information objects being preserved

Digital Migration Approaches

- Four primary types of digital migration in response to motivators, ordered by increasing risk of information loss:
 - Refreshment
 - Media replacement with no bit changes
 - Replication
 - No change to Packaging Information or Content Information bits (e.g. copying to new file / new location)
 - Repackaging
 - Some bit changes in Packaging Information (e.g multiple files packaged in directory structure get copied to other carrier)
 - Transformation
 - Reversible: Bit changes in Content Information are reversible by an algorithm
 - Non-reversible: Bit changes in Content Information are not reversible by an algorithm

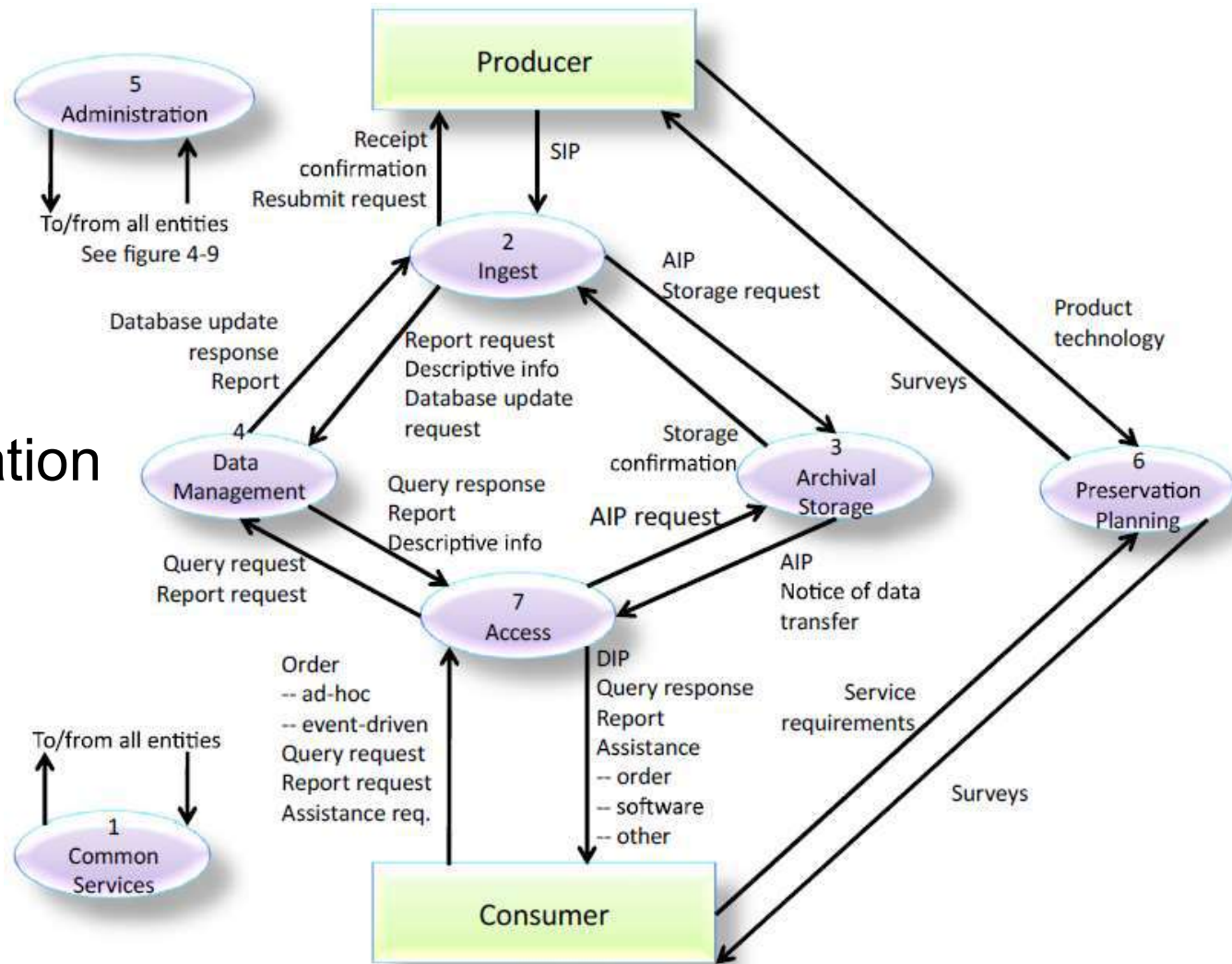


Digital Migration and AIPs

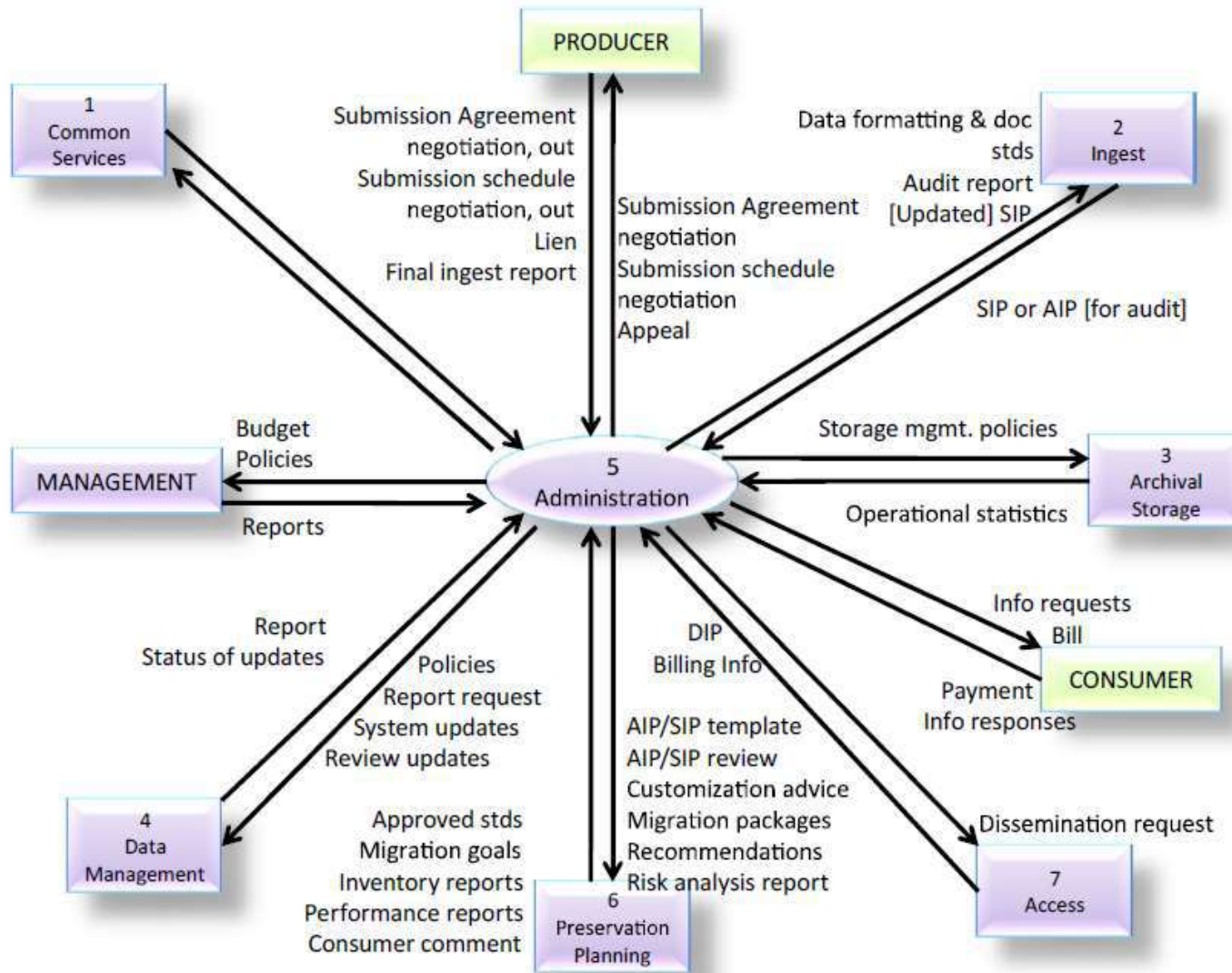
- Unless migration involves transformation:
no new AIP version
- Transformation:
new **AIP Version**
- Upgrading or improvement of AIPs: 
new **AIP Edition**
- Extracting or aggregating from multiple AIPs:
Derived AIP

OAIS Data Flow

- w/o administration



OAIS Data Flow



Common Services

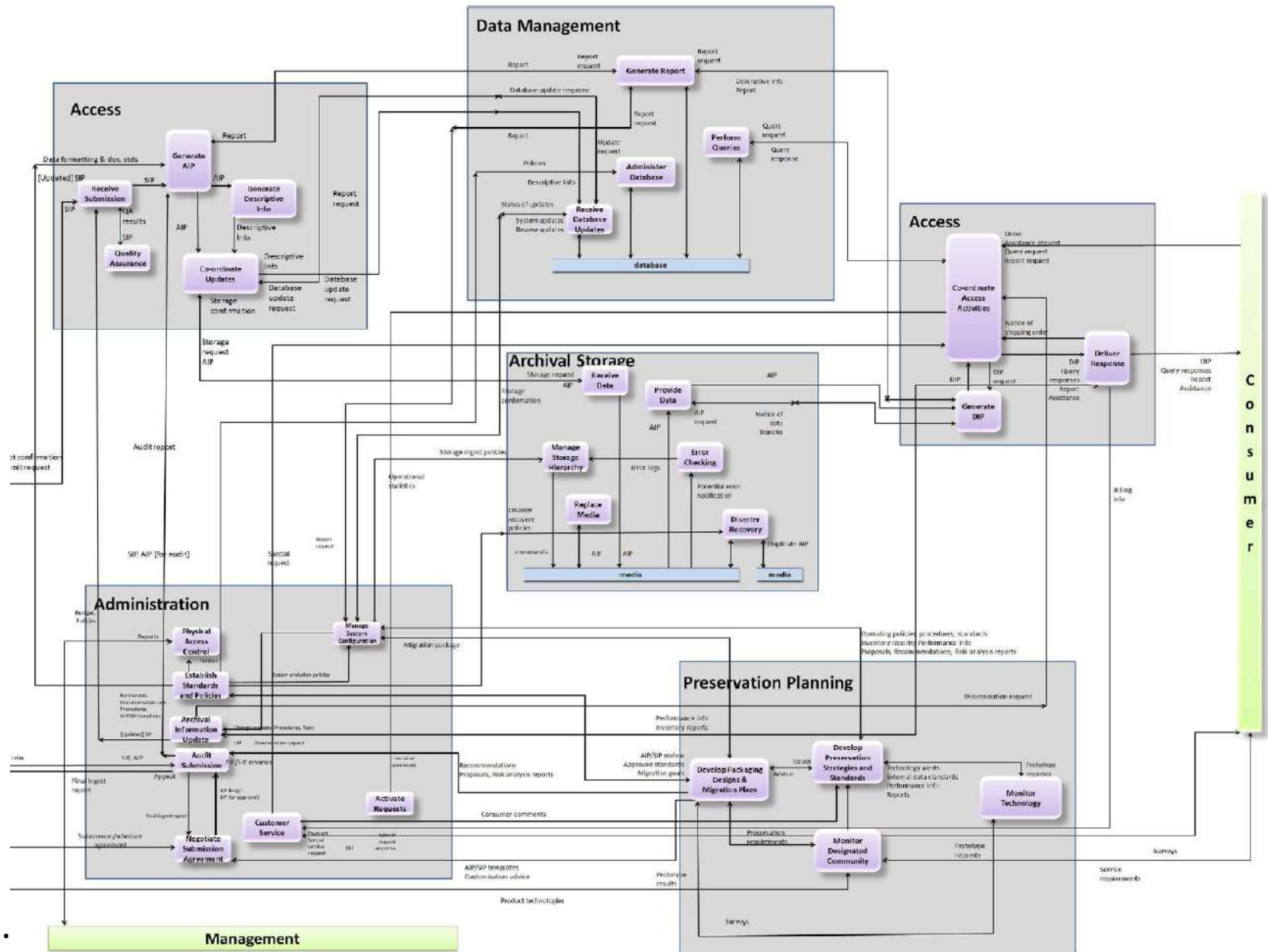
- Modern, distributed computing applications assume a number of supporting services
- Examples of Common Services include:
 - inter-process communication
 - name services
 - temporary storage allocation
 - exception handling
 - security
 - file and directory services

Common Services

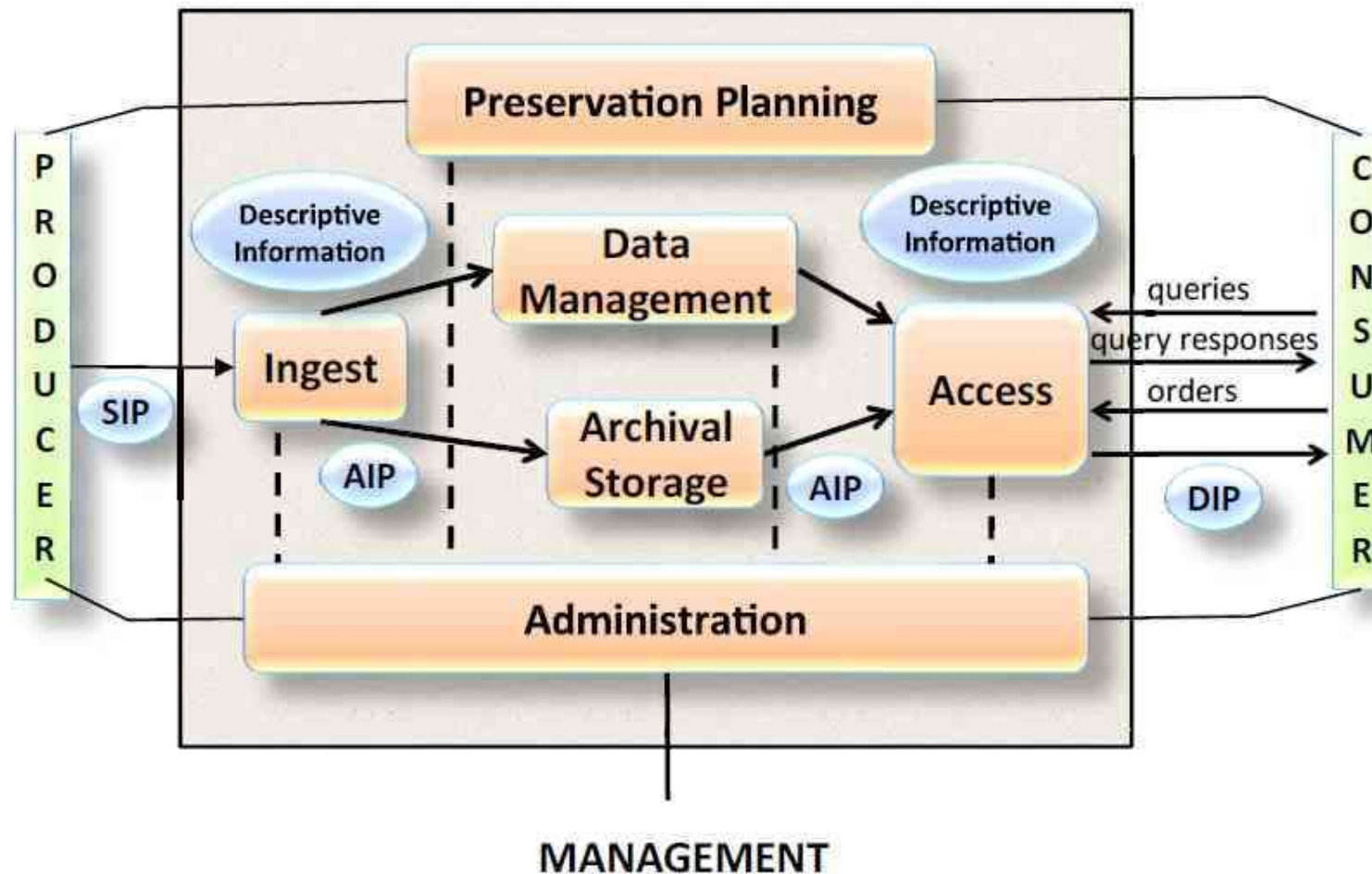
- Important:
 - All processing steps taken need to be documented (logs, protocols)
 - Reporting
 - Confirmations

 - This documentation is part of the archive as well

OAIS Composite Functional Entities



Open Archival Information System: Summary



SIP = Submission Information Package

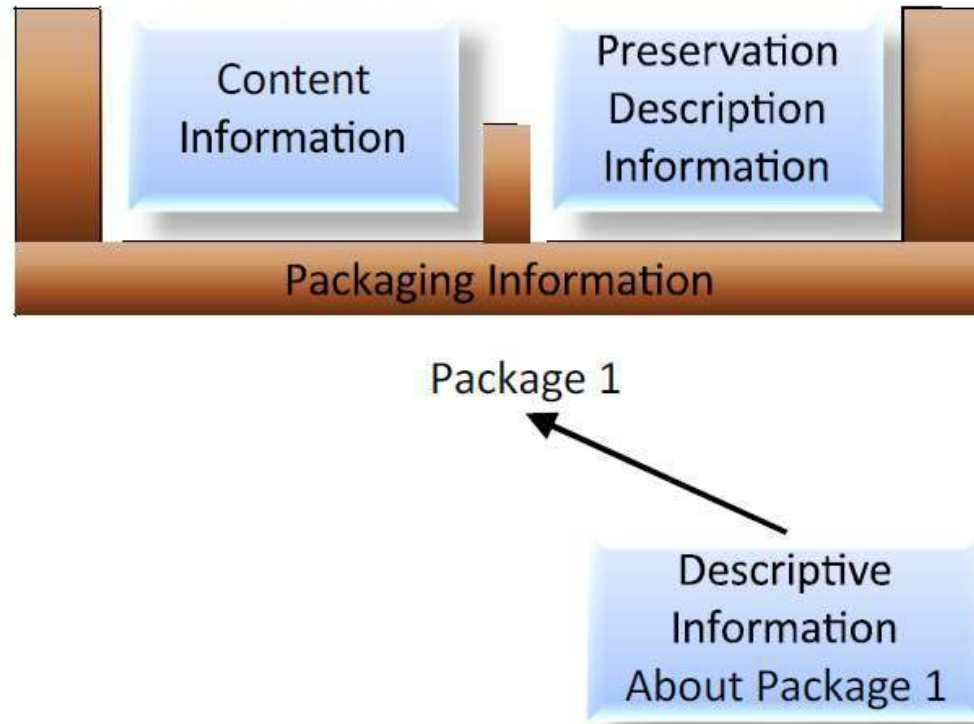
AIP = Archival Information Package

DIP = Dissemination Information Package

Outline

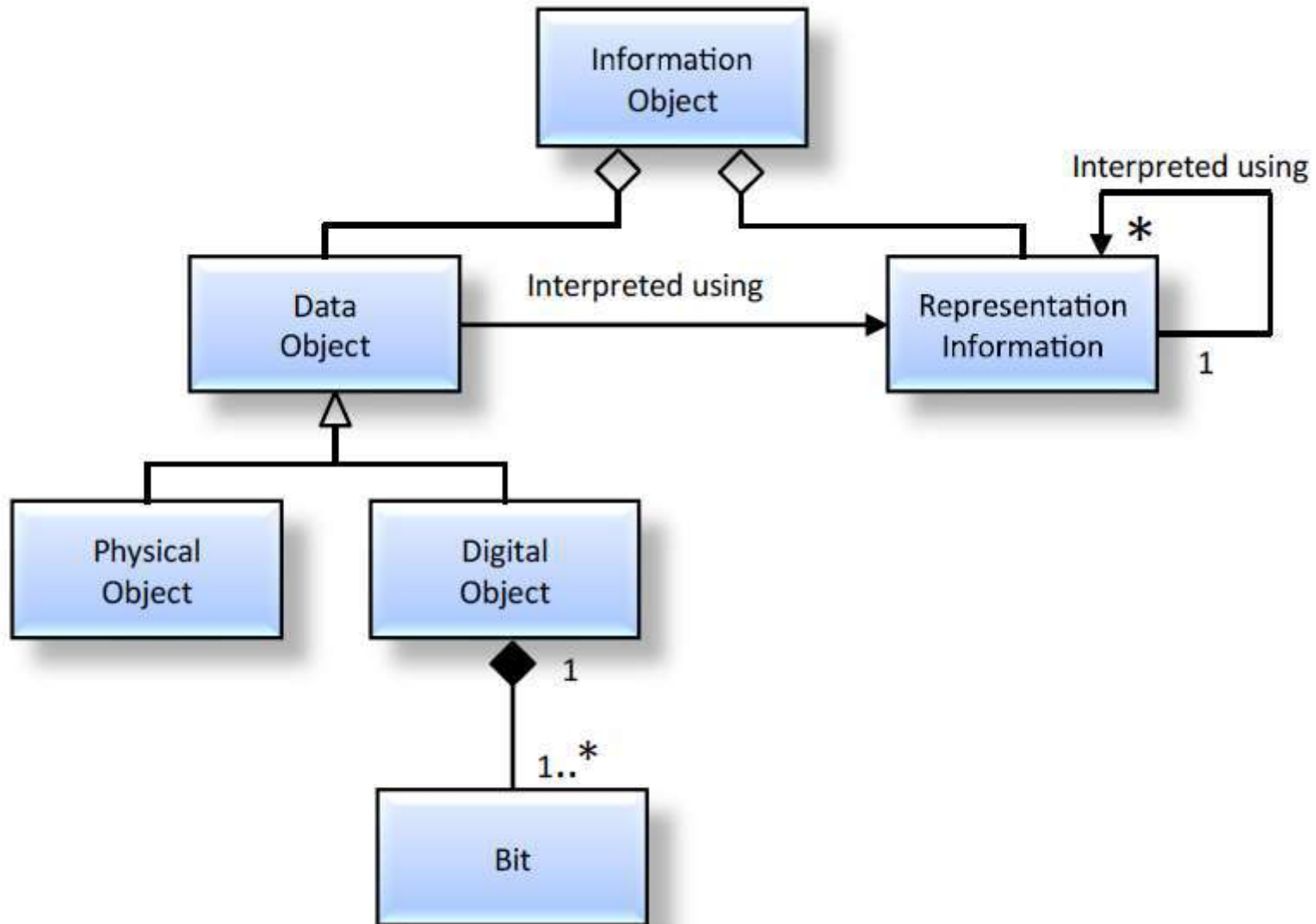
-
- Principles of the OAIS Model
 - Technical Overview
 - Functional Overview
 - **Information Modell**
 - Summary
-

Information Package Definition



- An Information Package is a conceptual container holding two types of information
 - Content Information
 - Preservation Description Information (PDI)
 - Plus descriptive information

Information Object



.....

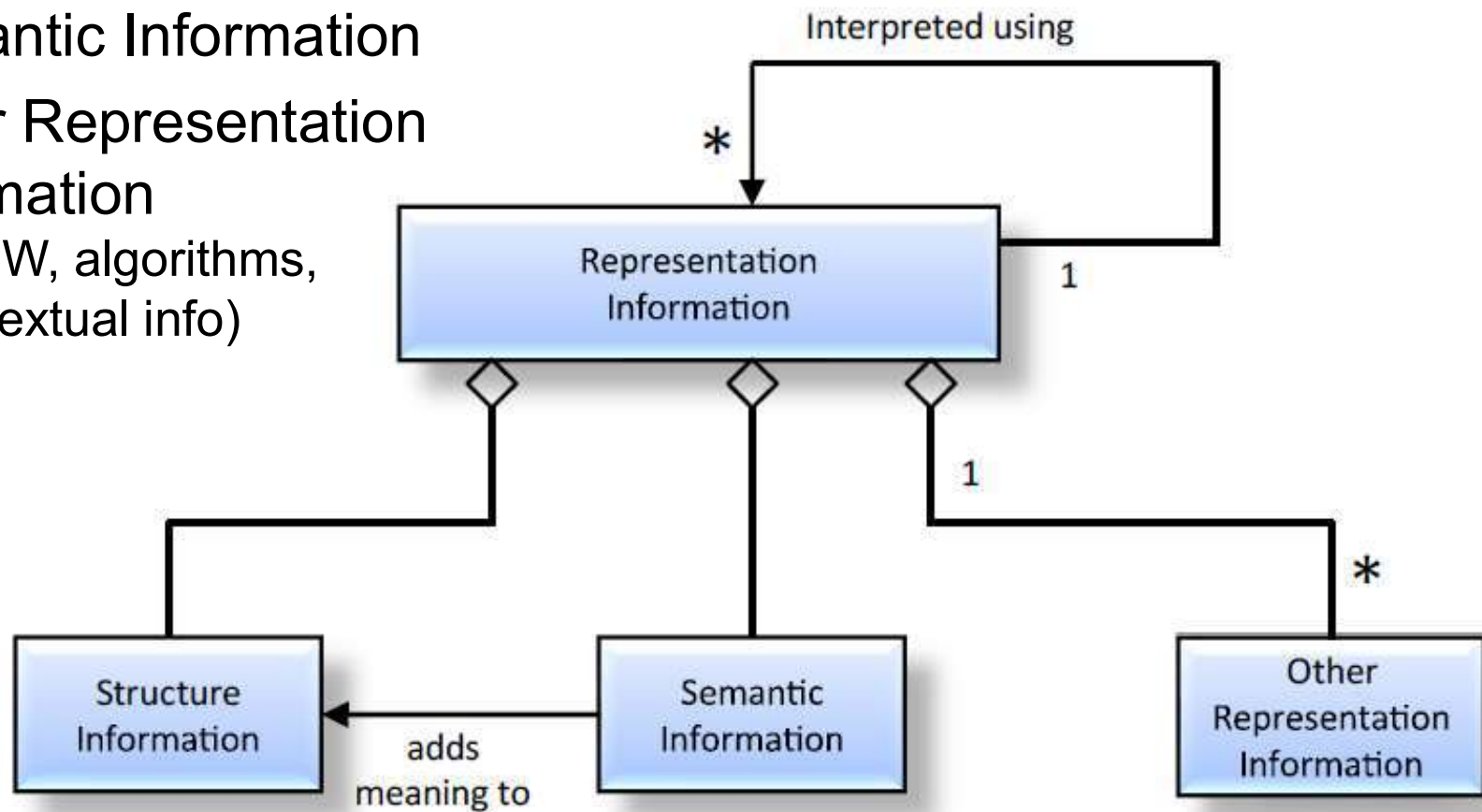
Representation Information

- The Representation Information accompanying a physical object, like a moon rock, may give additional meaning
 - It typically is a result of some analysis of the physically observable attributes of the rock
- The Representation Information accompanying a digital object, or sequence of bits, is used to provide additional meaning.
 - It typically maps the bits into commonly recognized data types such as character, integer, and real and into groups of these data types.
 - It associates these with higher level meanings which can have complex inter-relationships that are also described

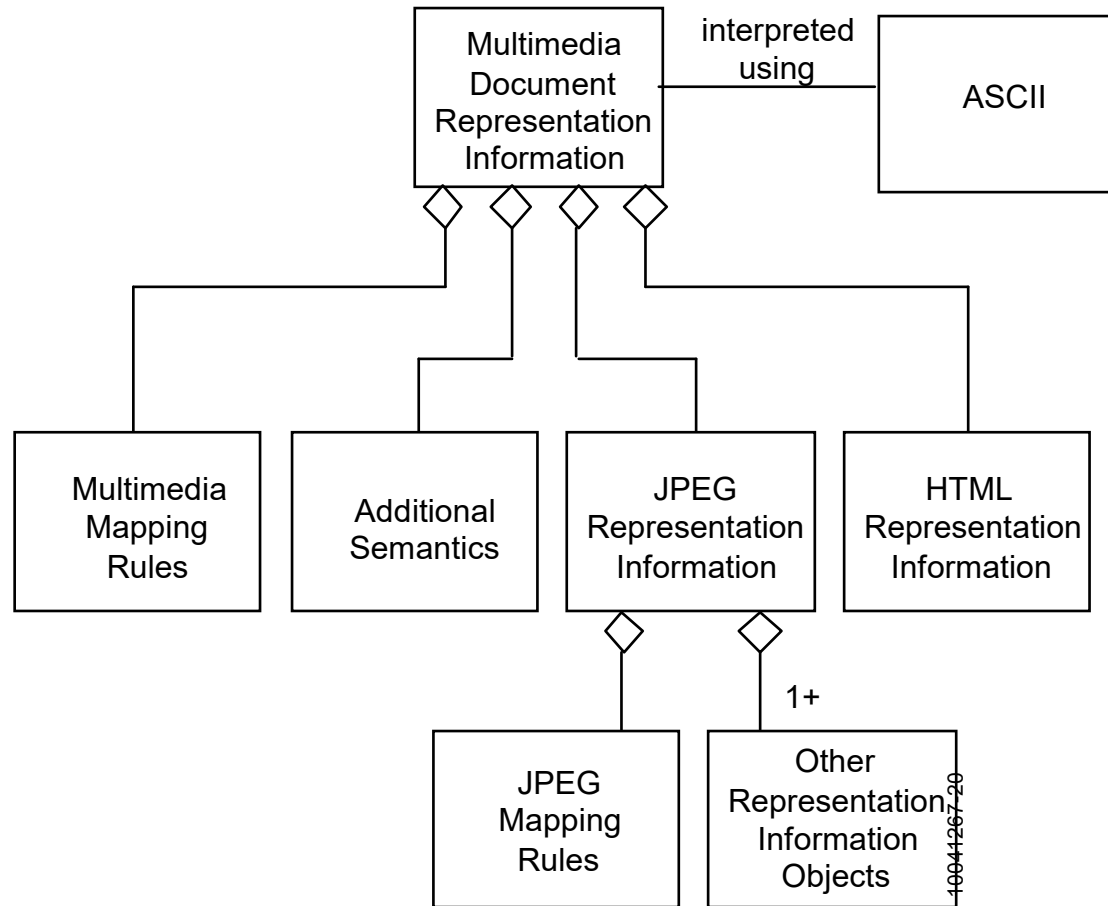
.....

Recursive Nature of Representation Information

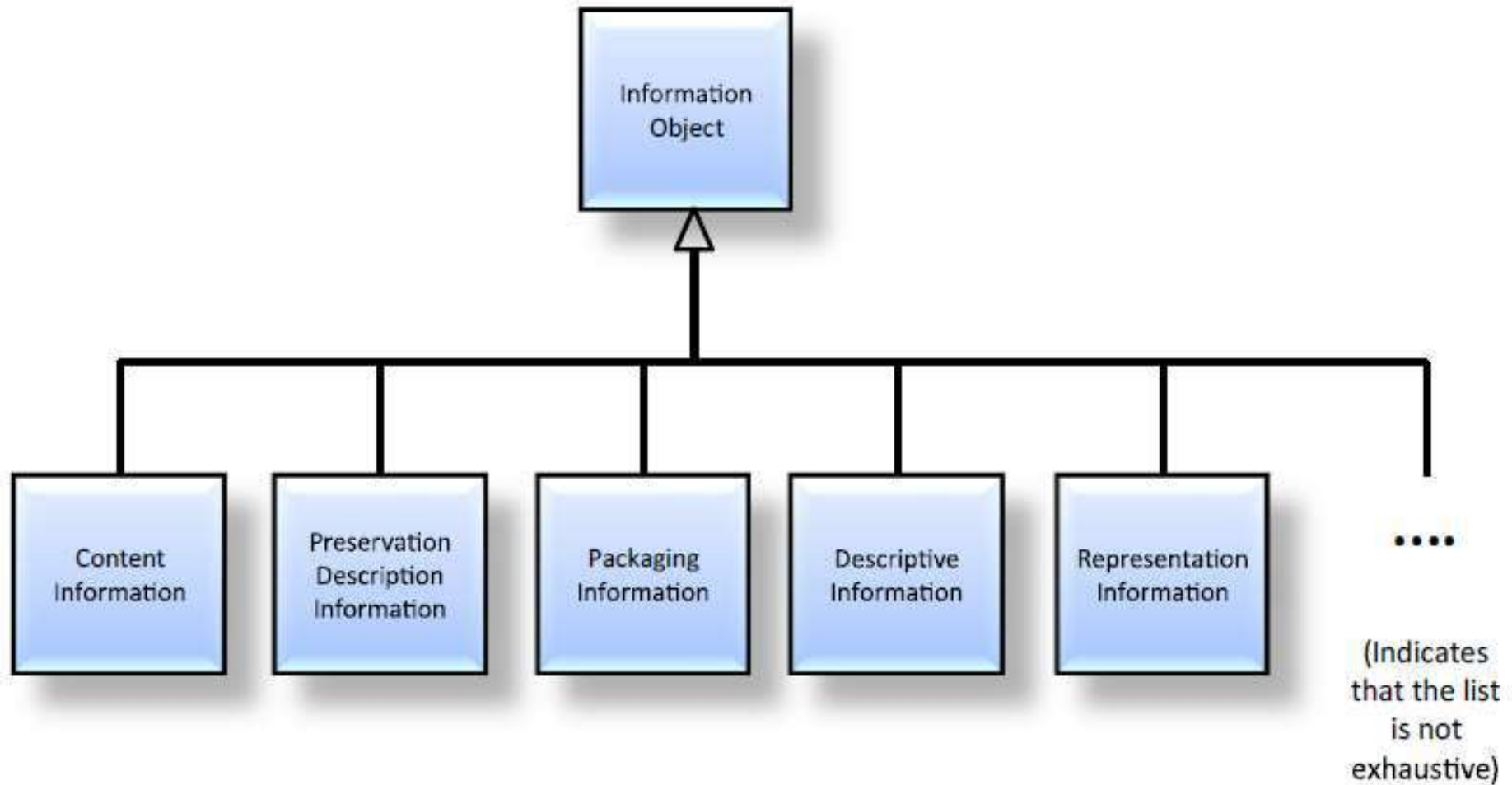
- Structure Information
- Semantic Information
- Other Representation Information
(e.g. SW, algorithms, other textual info)



Sample Representation Net



Types of Information Used in OAIIS



Content Information



- The information which is the primary object of preservation
- An instance of Content Information is the information that an archive is tasked to preserve.
- Deciding what is the Content Information may not be obvious and may need to be negotiated with the Producer
- The Data Object in the Content Information may be either a Digital Object or a Physical Object (e.g., a physical sample, microfilm)

- Provenance Information
 - Describes the source of Content Information, who has had custody of it, what is its history
- Context Information
 - Describes how the Content Information relates to other information outside the Information Package
- Reference Information
 - Provides one or more identifiers, or systems of identifiers, by which the Content Information may be uniquely identified
- Fixity Information
 - Protects the Content Information from undocumented alteration

PDI Examples

Content Information Type	Reference	Provenance	Context	Fixity	Access Rights
Space Science Data	<ul style="list-style-type: none"> Object identifier Journal reference Mission, instrument, title, attribute set 	<ul style="list-style-type: none"> Instrument description Principal Investigator Processing history Storage and handling history Sensor description Instrument Instrument mode Decommutation map Software interface specification Information Property Description 	<ul style="list-style-type: none"> Calibration history Related data sets Mission Funding history 	<ul style="list-style-type: none"> CRC Checksum Reed-Solomon coding 	<ul style="list-style-type: none"> Identification of the properly authorized Designated Community (Access Control) Permission grants for preservation and for distribution Pointers to Fixity and Provenance Information (e.g., digital signatures, and rights holders)

PDI Examples

<p>Digital Library Collections</p>	<ul style="list-style-type: none"> • Bibliographic description • Persistent identifier 	<ul style="list-style-type: none"> • For scanned collections: <ul style="list-style-type: none"> • metadata about the digitization process • pointer to master version • For born-digital publications: <ul style="list-style-type: none"> • pointer to the digital original • Metadata about the preservation process: <ul style="list-style-type: none"> • pointers to earlier versions of the collection item • change history • Information Property Description 	<ul style="list-style-type: none"> • Pointers to related documents in original environment at the time of publication 	<ul style="list-style-type: none"> • Digital signature • Checksum • Authenticity indicator 	<ul style="list-style-type: none"> • Legal framework(s) • Licensing offers • Specifications for rights enforcement measures applied at dissemination time • Permission grants for preservation and for distribution • Information about watermarking applied at submission and preservation time • Pointers to Fixity and Provenance Information (e.g., digital signatures, and rights holders)
---	--	--	--	---	---



PDI Examples

Content Information Type	Reference	Provenance	Context	Fixity	Access Rights
Software Package	<ul style="list-style-type: none"> Name Author/Originator Version number Serial number 	<ul style="list-style-type: none"> Revision history Registration Copyright Information Property Description 	<ul style="list-style-type: none"> Help file User guide Related software Language 	<ul style="list-style-type: none"> Certificate Checksum Encryption CRC 	<ul style="list-style-type: none"> Designated Community Legal framework(s) Licensing offers Specifications for rights enforcement measures applied at dissemination time Pointers to Fixity and Provenance Information (e.g., digital signatures, and rights holders)

Packaging Information



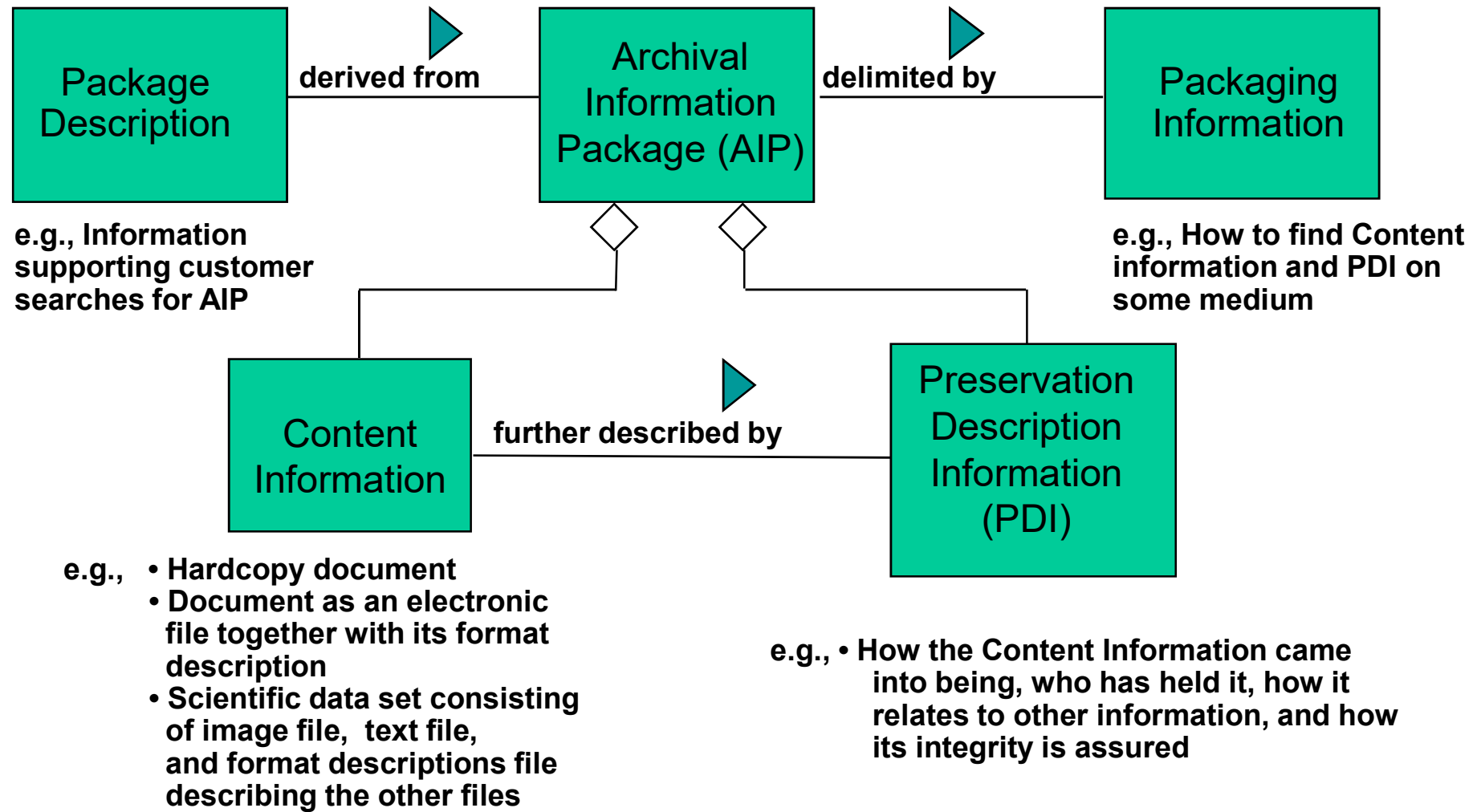
- Information which, either actually or logically, binds and relates the components of the package into an identifiable entity on specific media
- Examples of Packaging Information include tape marks, directory structures and filenames

Descriptive Information

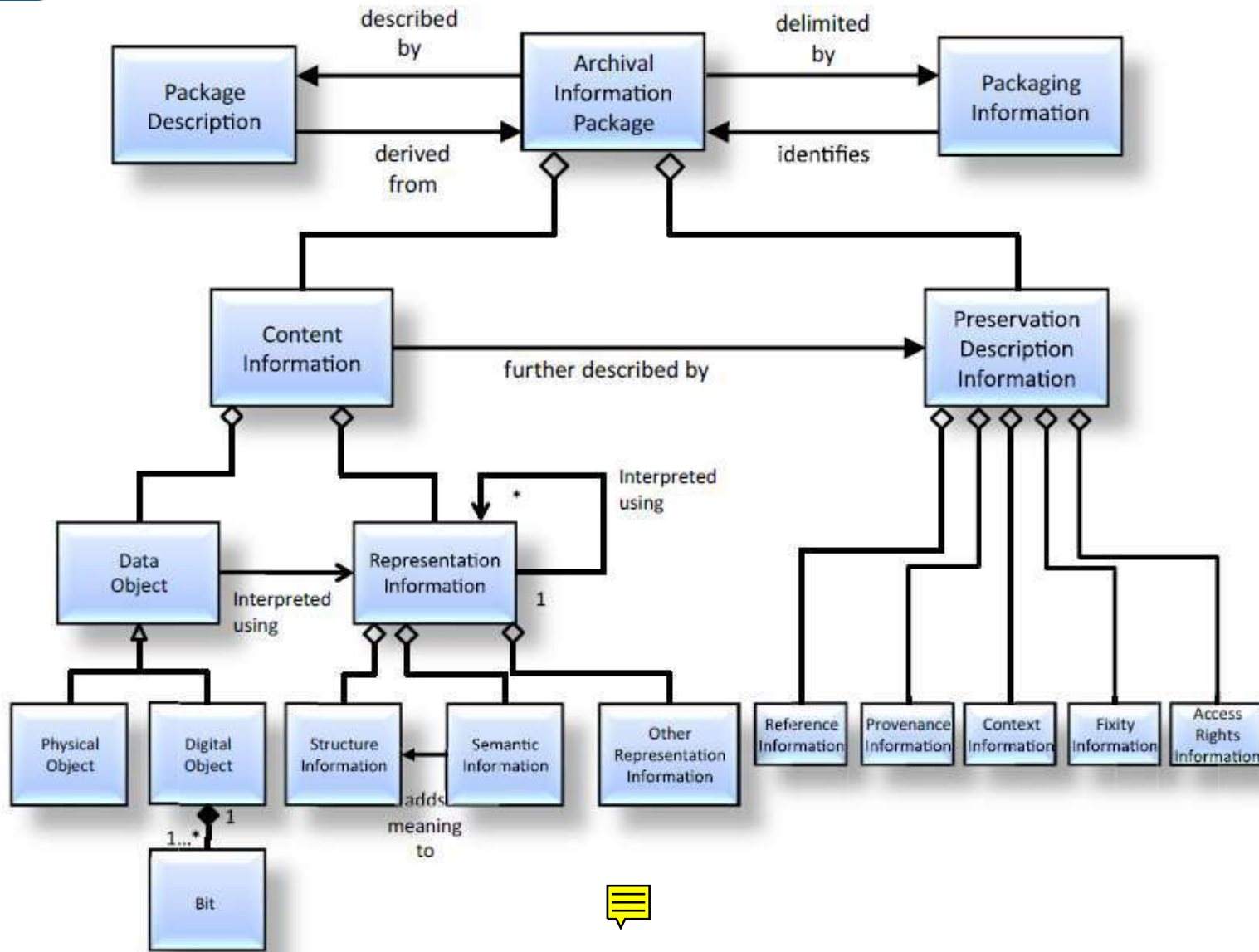


- Contain the data that serves as the input to documents or applications called Access Aids.
- Access Aids can be used by a consumer to locate, analyze, retrieve, or order information from the OAIS.

OAIS Archival Information Package

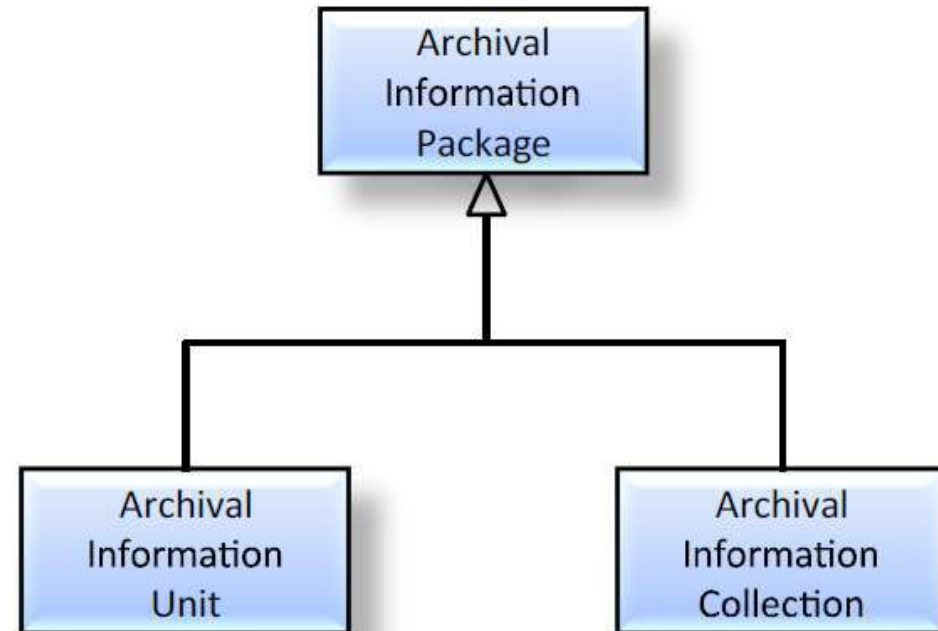


AIP detailed view



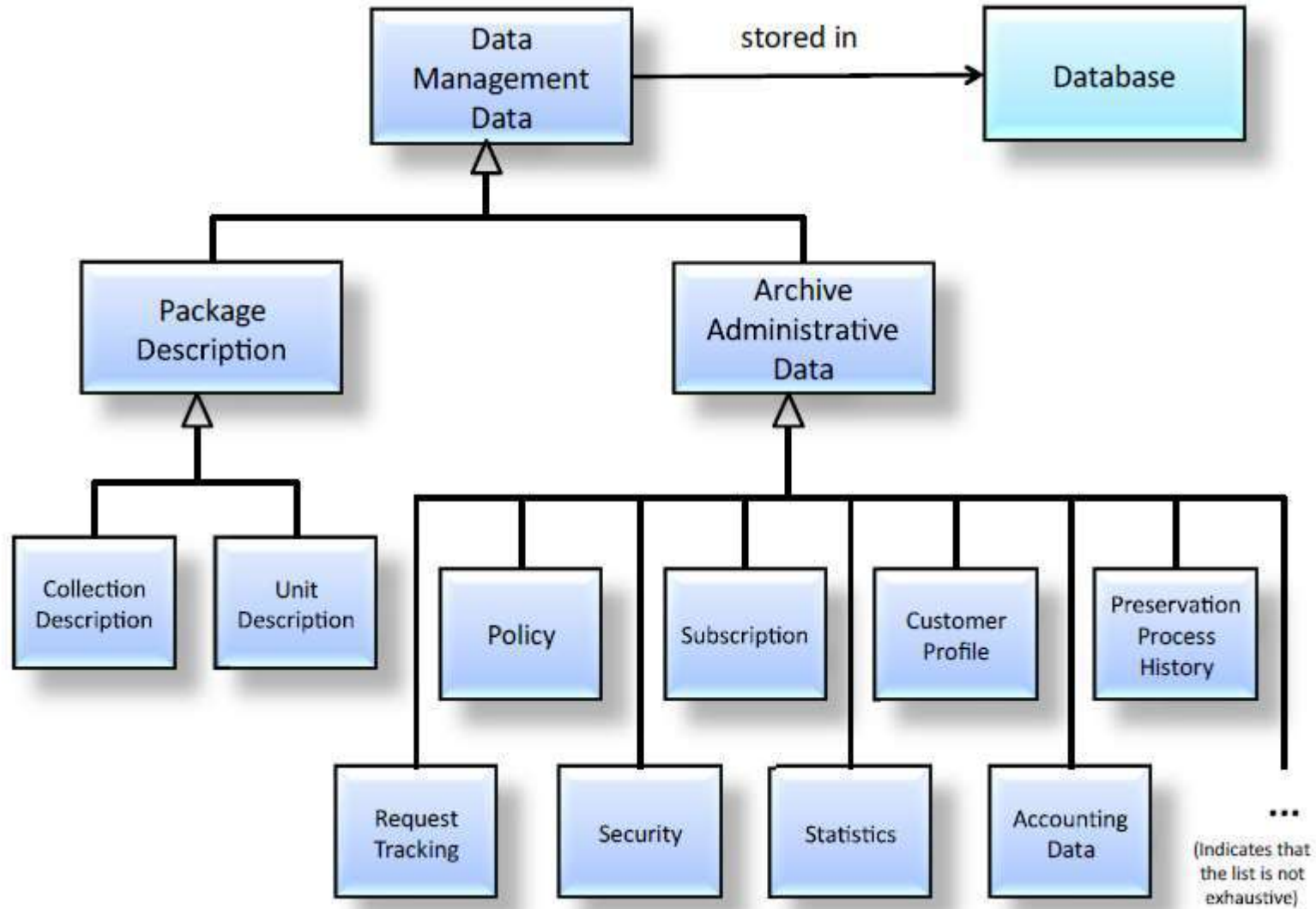
AIP Types

- Archival Information Unit (AIU) contains a single Data Object as the Content Object
- Archival Information Collection (AIC) contains multiple AIPs in its Content Object
 - Each member of an AIC is an AIP containing Content Information and PDI
 - The AIC contains unique PDI on the collection process



- Package Descriptions are needed by an OAIS to provide visibility and access to the OAIS holdings
- Package Descriptions contain 1 or more Associated Descriptions which describe the AIP Content Information from the point of view of a single Access Aid
- Some example of Access Aids Include:
 - Finding Aids - assist the consumer in locating information of interest
 - Ordering Aids - allow the consumer to discover the cost of and order AIUs of interest
 - Retrieval Aids - enable authorized users to retrieve the AIU described by the Unit Descriptor from Archival Storage

Data Management Information



.....

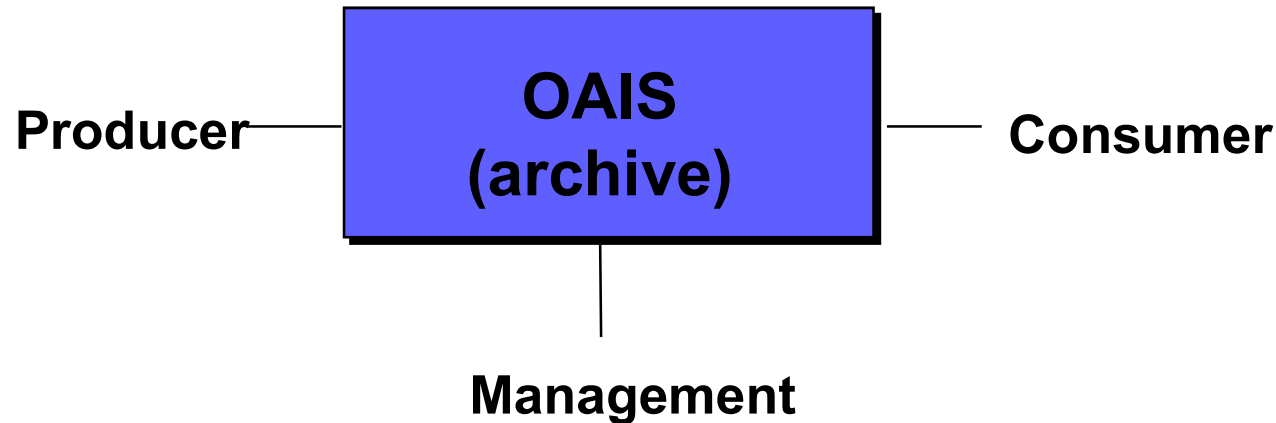
Information Model Summary

- Presented a model of information objects as containing data objects and representation objects
- Classified information required for Long-term archiving into 4 classes: Content Information, PDI, Packaging Information and Descriptive Information
- Described how these classes would be aggregated and related in an AIP to fully describe an instance of Content Information
- Presented information needed for Access, in addition to that needed for Long-term Preservation
- Put the Access oriented structures in the context of the other data needed to operate an OAIS

Outline

-
- Principles of the OAIS Model
 - Technical Overview
 - Functional Overview
 - Information Modell
 - Summary
-

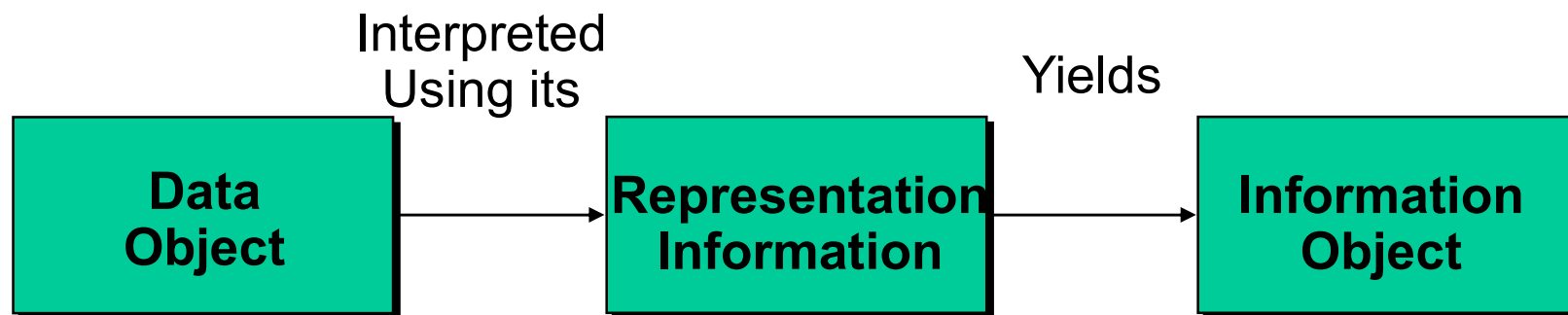
Model View of an OAIIS Environment



- Producer is the role played by those persons, or client systems, who provide the information to be preserved
- Management is the role played by those who set overall OAIIS policy as one component in a broader policy domain
- Consumer is the role played by those persons, or client systems, who interact with OAIIS services to find and acquire preserved information of interest

OAIS Information Definition

- Information is always expressed (i.e., represented) by some type of data
- Data interpreted using its Representation Information yields Information
- Information Object preservation requires clear identification and understanding of the Data Object and its associated Representation Information



Summary

- OAIS is a reference model
- OAIS no implementation specification
- Defines language, responsibilities, functionalities,...
- Can be used for all kind of archives, institutions, organizations, systems
- Can be used for all kinds of objects, physical or digital

Data Management Plans

Tomasz Miksa

tomasz.miksa@tuwien.ac.at

Agenda

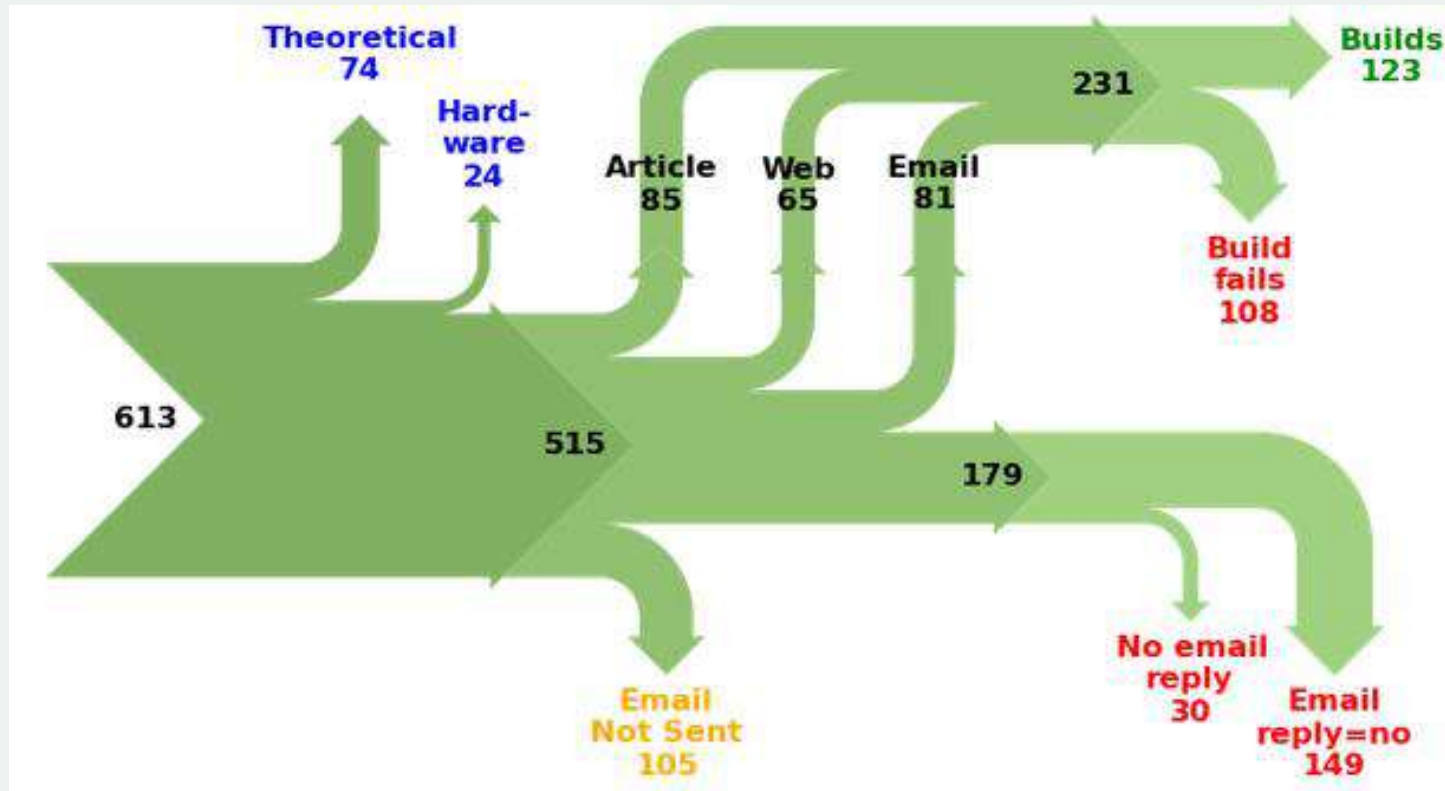
Why do we need to manage data properly?

What are Data Management Plans (DMPs)?

How to create a DMP?

Reproducibility - Computer Science

613 papers in 8 ACM conferences



Reproducibility - Computer Science

E-mail responses from authors

- Wrong version
- Code will be available soon
- Programmer left
- Bad backup practices
- Commercial code
- Proprietary academic code
- Intellectual property
- No intention to release
- ...

Variety of solutions

In response to these needs many solutions were proposed and are being implemented

- **FAIR principles**
- **open access** to scientific publications and data
- research **data repositories** to host the data
- **persistent identifiers** to locate the data
- **data management plans**
- ...

WHAT IS A DATA MANAGEMENT PLAN (DMP)?

Data Management Plan

DMP is a formal document

It outlines what you will do with your data **during** and **after** you complete your research

It ensures your data is safe for the **present** and the **future**

[from University of Virginia Library]




DMP is an awareness tool!

DMP makes you think

- what data you will use and where you get it from
- what infrastructure, software, licenses are needed
- what will be the output of your research
- how you will share your research outputs

DMP helps you organise yourself better

DMP can reveal how solid your methodology is

- is it a 'fishing expedition'? 




DMPs are used worldwide

- Required by
 - research funders
 - institutions, e.g. universities



Research Data Management

General Information

Research data management is an integral part of good research practice (see » [Research Integrity & Research Ethics](#)). The FWF therefore requires a data management plan (DMP) for all projects approved as of 1 January 2019. A DMP describes how data and their metadata are collected, organised, stored, published, shared, and archived for a specific project. Furthermore, the DMP outlines how the data will be made  FAIR, which means Findable, Accessible, Interoperable and Reusable. The » [FWF's Open Access Policy to Research Data](#) must be taken into account when drafting the DMP.

<https://www.fwf.ac.at/en/research-funding/open-access-policy/research-data-management>



Example: Projects funded by European Commission

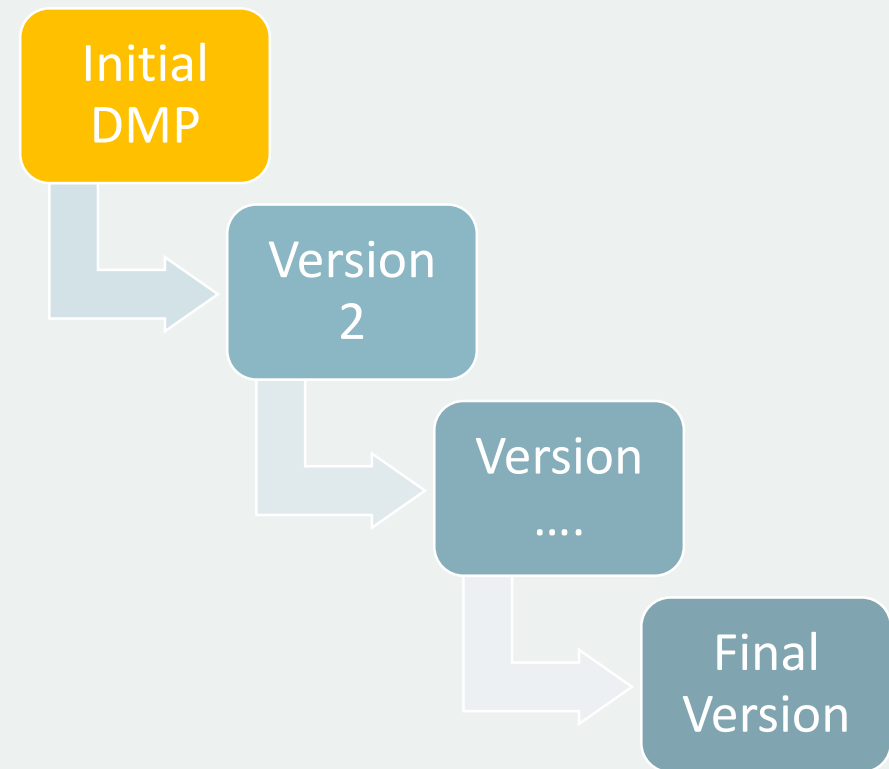
DMP is a **living document**

First version


- within the first 6 months

Updated versions

- when significant changes occur
 - new datasets
 - changes in policies
- periodic reporting
 - project reviews
- end of project



DMPs are not for research only

- DMPs are the requirement of research funders 
 - But it does not mean they are not useful elsewhere!
- Fire extinguisher is an obligatory car equipment
 - But it does not mean it is not useful elsewhere!
- Goal of this lecture
 - Not to make you experts on funder requirements
 - Examples provided are meant for illustration only
 - To help you create DMPs whenever you work with data
 - To improve your (Research) Data Management (R)DM!
- Overlaps with FAIR

HOW TO CREATE A DMP?

How to create a DMP?

Most cases by

- filling out a template
- answering questions from a checklist

Using software tools

- users choose appropriate funders template
- only relevant questions and guidance is presented

Science Europe Guidelines

Basis for many funder templates



Table of Contents	
Foreword by Dr Thierry Damerval	2
Introduction	4
GUIDANCE FOR ORGANISATIONS: CORE REQUIREMENTS FOR DATA MANAGEMENT PLANS	7
GUIDANCE FOR ORGANISATIONS: CRITERIA FOR THE SELECTION OF TRUSTWORTHY REPOSITORIES	11
GUIDANCE FOR RESEARCHERS: Translating the Core Requirements into a DMP template Guiding the Selection of Trustworthy Repositories	15
GUIDANCE FOR REVIEWERS: Evaluation Rubric for Data Management Plans	31
Notes and References	51
Annex: Compatibility with the FAIR Data Principles	52

4 For procedural elements of implementing DMPs, see the RDA DMP Common Standards Working Group: <https://www.rd-alliance.org/groups/dmp-common-standards-wg>





FWF Example

Based on the **SE requirements**

I General Information		
I.1 Administrative information	Provide information such as name of principal investigator, FWF project number, and version of DMP	<ul style="list-style-type: none"> - Provide the relevant grant information. - Consider regular updates of the DMP.
I.2 Data management responsibilities and resources	<p>Who (for example, role, position, and institution) will be responsible for data management?</p> <p>What resources will be dedicated to data management and ensuring that data will be FAIR (Findable, Accessible, Interoperable, Re-usable)?</p>	<ul style="list-style-type: none"> - Indicate who is responsible for implementing the DMP, and for ensuring it is reviewed and, if necessary, revised. - For collaborative projects, explain the co-ordination of data management responsibilities across partners. - Explain how the necessary resources (for example, time) to prepare the data for sharing/preservation have been costed in. Carefully consider and justify any resources needed to deliver the data. These may include storage costs, hardware, staff time, and repository charges.
II Data Characteristics		
II.1 Data description and collection or re-use of existing data	<p>How will new data be collected or produced and/or how will existing data be re-used?</p> <p>What data (types, formats, and volumes) will be collected or produced?</p>	<ul style="list-style-type: none"> - Explain which methodologies or software will be used if new data are collected or produced. - State any constraints on re-use of existing data if there are any. - Explain how data provenance will be documented. - Give details on the kind of data: for example, numeric (databases), textual (documents), image, audio, or video. - Give details on the data format: the way in which the data is encoded for storage, often reflected by the filename extension (for example, pdf, xls, doc, txt, or rdf).

DMP tools



DMP Online

- <https://dmponline.dcc.ac.uk/>

Data Stewardship Wizard

- <https://ds-wizard.org>

Argos

- <https://argos.openaire.eu/splash/>

RDMO

- <https://rdmorganiser.github.io/en/>

A screenshot of the DMP ONLINE web application. The interface is primarily orange and white. At the top, there's a navigation bar with 'DMP ONLINE' logo and links for 'View plans', 'Create plan', 'About', 'Future plans', 'Help', and 'Change language'. Below this, the main content area shows a quiz titled 'FFG Webinar Horizon 2020 Example' with a progress indicator '0/71 questions answered'. The quiz is divided into sections: '1. Data summary (1 question, 0 answered)', '2. FAIR data (4 questions, 0 answered)', and '2.1 Making data findable, including provisions for metadata:'. Under section 2.1, there are several bullet points: 'Outline the discoverability of data (metadata provision)', 'Outline the identifiability of data and refer to standard identification mechanism. Do you make use of persistent and unique identifiers such as Digital Object Identifiers?', 'Outline naming conventions used', 'Outline the approach towards search keyword', 'Outline the approach for clear versioning', and 'Specify standards for metadata creation (if any). If there are no standards in your discipline describe what metadata will be created and how'. Below the text is a rich text editor with a toolbar containing bold, italic, list, link, and image icons. A 'Save' button is located at the bottom left of the editor area. On the right side, there's a 'Guidance' panel titled 'EC Guidance' which contains a note: 'The Research Data Alliance provides a Metadata Standards Directory that can be searched for discipline-specific standards and associated tools.'

DMP tool at TUW

The screenshot shows the DAMAP web application interface. At the top left, there is a blue header with 'DAMAP' and 'EN'. Below it, a dark grey bar displays the user's name 'Tomasz Miksa' and a 'Logout' link. The top right features navigation links for 'HOME' and 'PLANS', along with the 'TU WIEN' logo. The main content area is titled 'Data Management Plan' and includes a yellow warning box for the 'Test environment'. Below this, the 'DAMAP' section provides a welcome message and two buttons: 'My plans' and '+ Create new plan'. A 'What is a DMP?' section explains the purpose of a DMP, and a 'TU DMP Tool' section lists its features and capabilities.

DAMAP EN

HOME PLANS TU WIEN

Tomasz Miksa

Logout

Data Management Plan

Test environment

This application instance is for development and testing purposes only. Content may be deleted at any point in time without prior notification.

DAMAP

Welcome to TU DMP Tool, a service that helps you to create and update the Data Management Plan (DMP) for your project.

[My plans](#) [+ Create new plan](#)

What is a DMP?

A [DMP](#) is a structured document that keeps record of what research data is created and what happens to that data during and after a project. It helps with planning the research process, managing your data in accordance with the [FAIR Principles](#), and defining rights and responsibilities in a research project involving several researchers or institutions.

TU DMP Tool

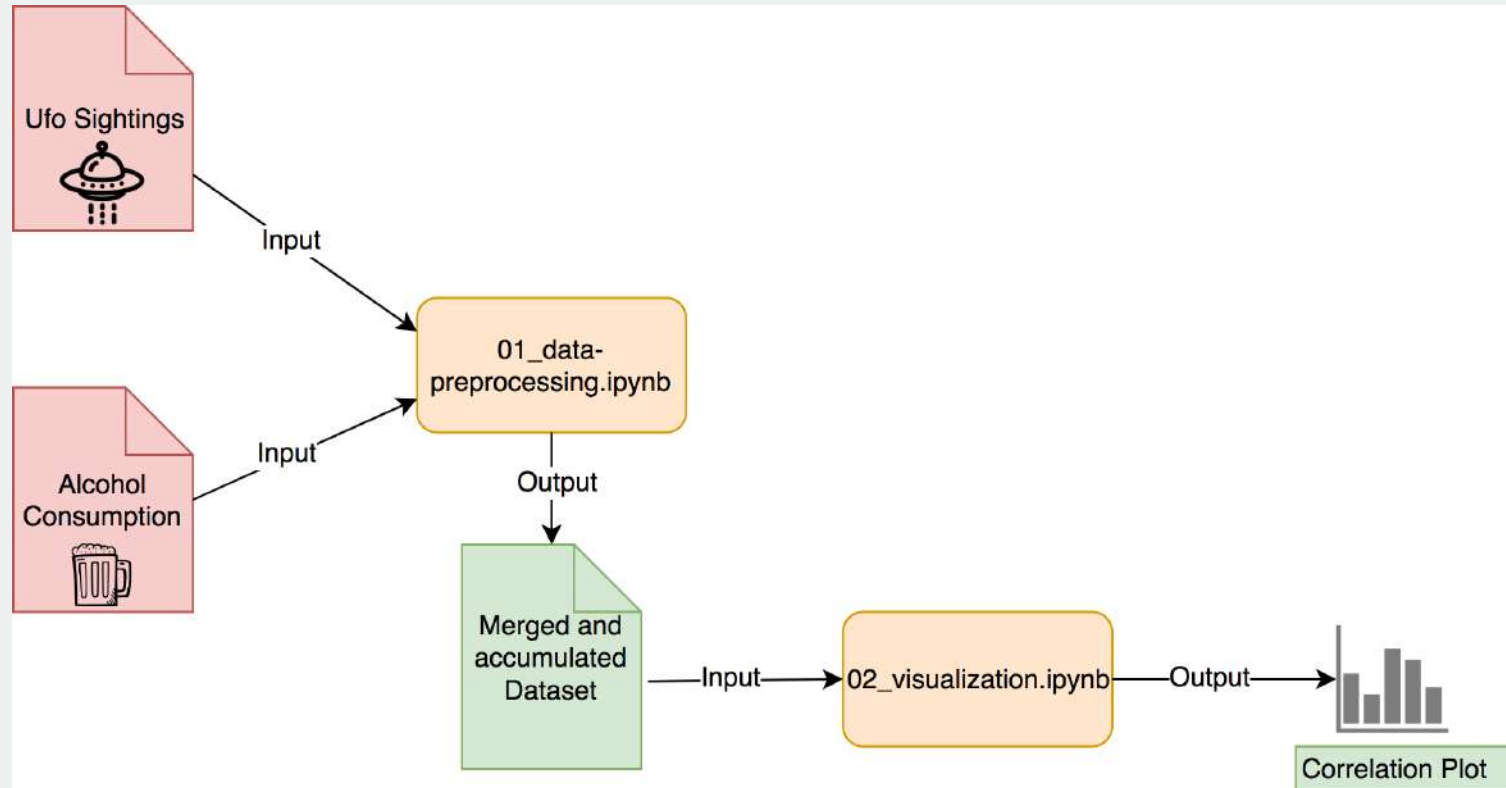
- guides you step by step through the different sections of a DMP following the [Science Europe Practical Guide](#)
- exports a pre-filled DMP as a Word document that you can customize and use for submission to European and national funders, for example FWF and FFG
- saves you work by
 - pre-filling content with detailed information from TISS and other systems
 - providing wizards, guidance, and item lists to choose from
 - suggesting answers that you can either comply with or adjust to your needs
- is compatible with the [RDA recommendation on machine actionable DMPs](#).

Version: 1.1.0

(and what I should also do!)

WHAT SHOULD I WRITE IN FACT?

Correlating Alcohol Consumption and UFO Sightings in the USA (running example)



Running example

This example is ***very simplified***

- Few inputs, few outputs
 - Compared to thousands of files processed
- ‘Digital objects’ are identical with files (which is not the case usually)
 - Compared to collections of files with a complex structure
- Experiment is identical with a single notebook
 - Compared to several workflows, software tools, etc.

DP Exercise

- Do not use the example as an easy copy-paste!

FWF Template – running example

I. General information

II. Data Characteristics

III. Documentation and Data Quality

IV. Data Storage, Sharing and Long-Term Preservation

V. Legal and Ethical Aspects

This lecture: not an exhaustive walk-through! Only interesting/relevant aspects!

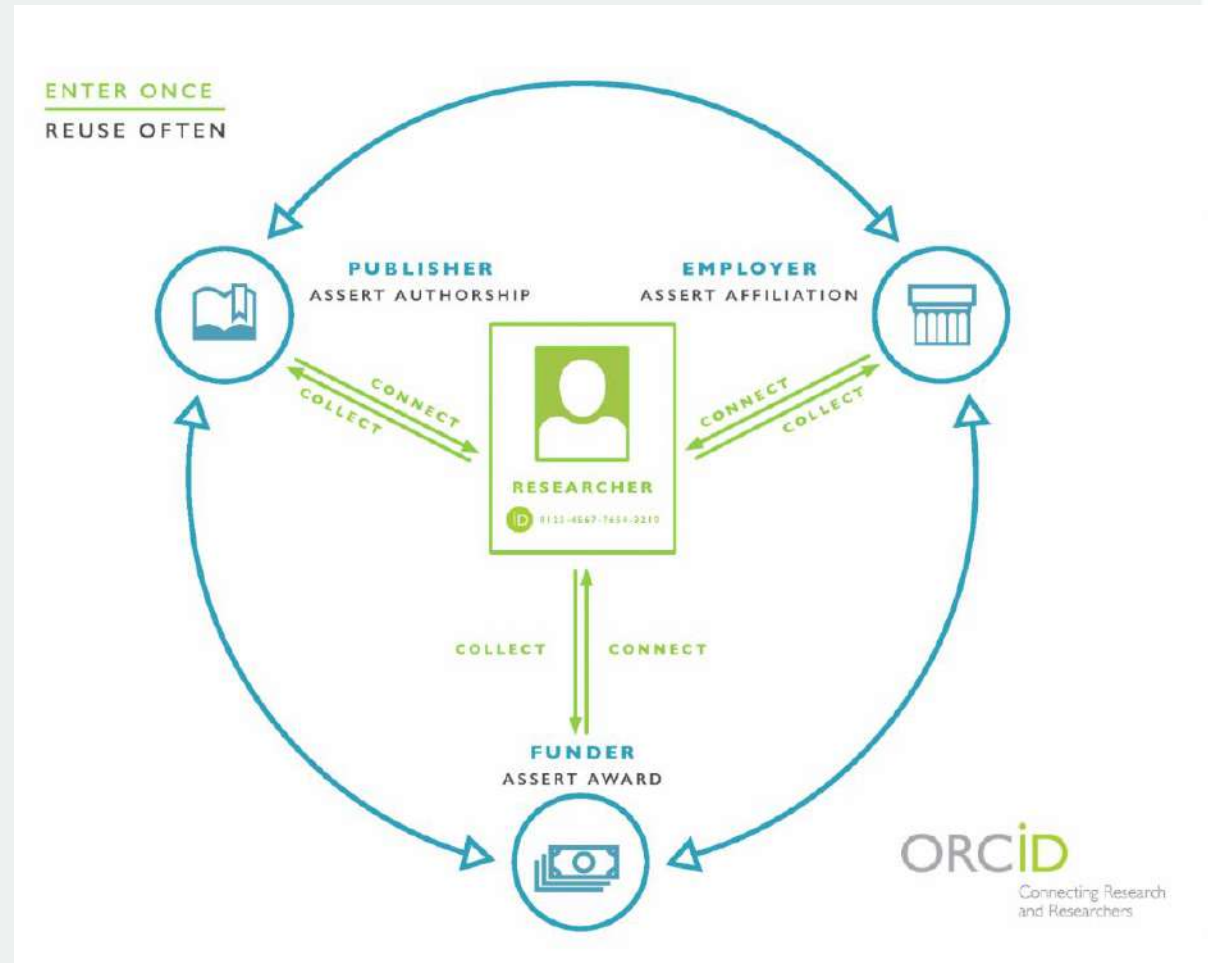
I General Information	
I.1 Administrative information	
I.2 Data management responsibilities and resources	
II Data Characteristics	
II.1 Data description and collection or re-use of existing data	
III Documentation and Data Quality	
III.1 Metadata and documentation	
III.2 Data quality control	
IV Data Storage, Sharing, and Long-Term Preservation	
IV.1 Data storage and backup during the research process	
IV.2 Data sharing and long-term preservation	
V Legal and Ethical Aspects	
V.1 Legal aspects	
V.2 Ethical aspects	

I. GENERAL INFORMATION

ORCID – persistent identifier for people

ORCID ID

- Unique person ID
- ORCID assigned once
- Person can change affiliations (jobs)
- Example: 0000-0002-4929-7875



Search English

ORCID Connecting Research and Researchers

4,115,029 ORCID iDs and counting. [See more...](#)

Daniel Mietchen

ORCID ID
<https://orcid.org/0000-0001-9488-1870>

[Print view](#)

Also known as
 D. Mietchen, Mietchen, Daniel, Mietchen, D., EvoMRI, D Mietchen, Mietchen D, Mietchen-D

Country
 Germany

Keywords
 open science, open data, open access, magnetic resonance microscopy, evolution, biodiversity, social machines, vocal learning

Websites
[Twitter](#)
[Wikidata](#), [Wikipedia et al.](#)
[GitHub](#)
[Open Science Q & A](#)
[Scholia](#)

Other IDs
 Scopus Author ID: 7801384320
 ResearcherID: A-7748-2009

Employment (2) Sort

National Center for Biotechnology Information: Bethesda, MD, United States
 2015-03-01 to present | Intramural researcher (Computational Biology Branch)
 Source: Daniel Mietchen

Museum für Naturkunde - Leibniz-Institut für Evolutions- und Biodiversitätsforschung: Berlin, Berlin, Germany
 2013-08-16 to 2015-02-28 | Researcher (Digital World)
 Source: Daniel Mietchen

Works (64) Sort

Machine-actionable data management plans (maDMPs)
 Research Ideas and Outcomes
 2017-04-05 | journal-article
 DOI: [10.3897/rio.3.e13086](https://doi.org/10.3897/rio.3.e13086)
 Source: CrossRef Metadata Search Preferred source

Progress in promoting data sharing in public health emergencies
 Bulletin of the World Health Organization
 2017-04-01 | journal-article
 DOI: [10.2471/blt.17.192096](https://doi.org/10.2471/blt.17.192096)
 Source: CrossRef Metadata Search Preferred source


Strategies and guidelines for scholarly publishing of biodiversity data

II. DATA CHARACTERISTICS

slido

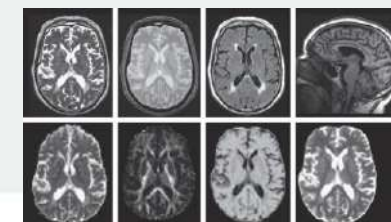
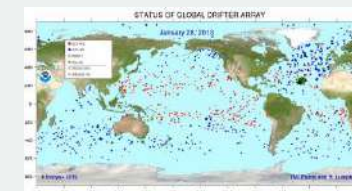
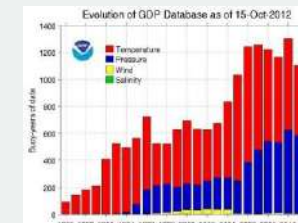
What is data?



 Start presenting to display the poll results on this slide.

What is data?

- Instrument measurements
- **Experimental observations**
- Still images, video and audio
- Text documents, spreadsheets, databases
- Quantitative data (e.g. survey data)
- Survey results & interview transcripts
- Simulation data, models & software
- Slides, artefacts, specimens, samples
- Questionnaires
- Sketches, diaries, lab notebooks ...



Data Summary

Type

- text, spreadsheets, software, models, images, movies, audio, patient records, etc.

Source

- human observation, laboratory, field instruments, experiments, simulations, compilations, etc.

Volume

- total volume of data, number of files, etc.

Data and file formats

- non-proprietary formats
- used within community



Data Summary - example

Produced Data

This project produces aggregated dataset in CSV format (Filesize ~800K) that contains data points that combine alcohol consumption data with the UFO sighting data and a correlation plot of these in PNG format (Filesize (~100K)).

Input Data

Project accesses two external CSV datasets. Both datasets have been downloaded and saved along with the source code in the folder *data/raw*.

1. Alcohol Consumption: OECD (2018), Alcohol consumption (indicator).

DOI: 10.1787/e6895909-en (Accessed on 22 March 2018)

File Location: `data/raw/DP_LIVE_22032018202902423.csv`

File Size: 112K



Data Summary - example

The experiment has been conducted with Jupyter notebooks. The notebooks contain the experiment's code, accompanying documentation, tables and plots.

We have included instructions (README.md) on how to run the experiment either directly or via Docker.

Running the code

To run the code in this repository you will need to have access to a machine running `python` (at least version 3.5) and `pip`.

Run `pip install -r requirements.txt` to install the required dependencies.

Once the dependencies have been installed, start the jupyter notebook server via `jupyter notebook` and open <http://localhost:8888>.

In the `notebooks` folder you'll find the following notebooks:

`01_data-preprocessing.ipynb`

Running this notebook generates a dataset consisting of the number of ufo sightings and the alcohol consumption in the usa per year by preprocessing and accumulating the data provided by the datasources mentioned above.

III. DOCUMENTATION AND DATA QUALITY

What is in the picture?



slido



How is this bird called?

ⓘ Start presenting to display the poll results on this slide.

slido



What is the sex of this bird?

ⓘ Start presenting to display the poll results on this slide.

slido



On which continent was the picture taken?

ⓘ Start presenting to display the poll results on this slide.

Metadata – Atlas Of Living


Atlas Of Living Australia ALA Apps ALA Info Search the Atlas Search

NatureShare - 2380_Gymnorhina_tibicen
HumanObservation of *Cracticus tibicen* | Australian Magpie recorded on 2011-04-17T12:32:00+1000


Flag an Issue Contact curator

Dataset
Event
Taxonomy
Geospatial
Images
Data quality tests (1 4 21 13 48)
Additional political boundaries information
Environmental sampling for this location

Location of record



Images



Photographer: Russell Best

Dataset

Date resource	NatureShare
Catalogue number	2380_Gymnorhina_tibicen
Basis of record	Human observation
Observer	Best, R. Russel Supplied as 'Russell Best'
Rights	CC BY 2.5 AU
More details	http://natureshare.org.au/observation/2380/
Photographer	Russell Best
Rights holder	Russell Best, via NatureShare
Occurrence remarks	Tags: Female
Occurrence status	present
Abcd identification qualifier	Not provided

Event

Record date	[date not supplied] Supplied date: '2011-04-17T12:32:00+1000'
Event remarks	Photo date/time used

Taxonomy

Scientific name	<i>Cracticus tibicen</i> Supplied scientific name: 'Gymnorhina tibicen'
Taxon rank	Species
Common name	Australian Magpie
Kingdom	Animalia
Phylum	Chordata
Class	Aves
Order	Passeriformes
Family	Artamidae
Genus	<i>Cracticus</i>
Species	<i>Cracticus tibicen</i>

Metadata – Atlas Of Living Australia

Dataset

Data resource	NatureShare
Catalogue number	2380_Gymnorhina_tibicen
Basis of record	Human observation
Observer	Best, R. Russell <i>Supplied as "Russell Best"</i>
Rights	CC BY 2.5 AU
More details	http://natureshare.org.au/observation/2380/
Photographer	Russell Best
Rightsholder	Russell Best via NatureShare
Occurrence remarks	Tags: Female
Occurrence status	present
Abcd identification qualifier	Not provided

Metadata – Atlas Of Living Australia

Event

Record date	[date not supplied] <i>Supplied date "2011-04-17T12:32:00+1000"</i>
Event remarks	Photo date/time used.

Taxonomy

Scientific name	<i>Cracticus tibicen</i> <i>Supplied scientific name "Gymnorhina tibicen"</i>
Taxon rank	Species
Common name	Australian Magpie
Kingdom	Animalia
Phylum	Chordata
Class	Aves
Order	Passeriformes
Family	Artamidae
Genus	<i>Cracticus</i>
Species	<i>Cracticus tibicen</i>
Taxonomic issues	No issues
Name match metric	Exact match The supplied name matched the name exactly.

Metadata – Atlas Of Living Australia

Geospatial

Country	Australia
State or territory	Victoria
Local government area	Macedon Ranges (S)
Latitude	-37.421078
Longitude	144.61954
Geodetic datum	EPSG:4326
Biome	Terrestrial
Verbatim longitude	144.619541
Verbatim latitude	-37.421077

Location of record



Standards and metadata

Metadata

- helps to understand and interpret data
- provides details about experiment setup
 - who, when, in which conditions, tools, versions, etc.
- helps identify and discover new data

Use community standards to enable interoperability

Metadata is also covered in the lecture on FAIR

<http://www.dcc.ac.uk/resources/metadata-standards>



Metadata - example

The metadata file can be found inside the project folder (/documentation/metadata.xml).

- experiment title, authors, date, tools, coverage, rights, etc.

Additionally a descriptive file is added, which explains the axes and units used in the output files, this file can be found inside the project folder as well (/documentation/description.txt).

- The alcohol consumption is the average consumption rate in liters/capita of USA inhabitants, which are older than 17.

Metadata-example

```
<?xml version="1.0" encoding="UTF-8"?>

<metadata xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:dc="http://purl.org/dc/terms/">
  <dc:title>UFOs & Alcohol</dc:title>
  <dc:creator>Marc Dietrichstein (https://orcid.org/0000-0003-4890-3498)</dc:creator>
  <dc:creator>Markus Neumeyer (https://orcid.org/0000-0002-4081-0716)</dc:creator>
  <dc:subject>Correlation of alcohol consumption and UFO sightings</dc:subject>
  <dc:description>Automated tool that investigates and computes the correlation between UFO sightings and alcohol consumption</dc:description>
  <dc:date>23.03.2018</dc:date>
  <dc:type>DataGeneration</dc:type>
  <dc:format>Jupyternotebook</dc:format>
  <dc:source>Ufo Sightings: Sigmond Axel. (2014)</dc:source>
  <dc:source>Alcohol Consumption: OECD (2018)</dc:source>
  <dc:language>English</dc:language>
  <dc:coverage>1960 - 2014 </dc:coverage>
  <dc:rights>Free access</dc:rights>
</metadata>
```

Some comments

- Sometimes it is hard to tell the difference between data and metadata
 - [NetCDF](#) (network Common Data Form)
 - is a file format for storing multidimensional scientific data (variables) such as temperature, humidity, pressure, wind speed, and direction.
 - Is self-describing, meaning that a netCDF file includes information about the data it contains, such as when data elements were captured and what units of measurement were used.
 - https://youtu.be/K1_8EqCJlwo
- Very often metadata is reduced to domain independent metadata
 - e.g. only Dublin Core or Data Cite metadata in a repository
 - Author, title, description, license, etc.
 - What about domain specific metadata?

IV. DATA STORAGE, SHARING AND LONG-TERM PRESERVATION

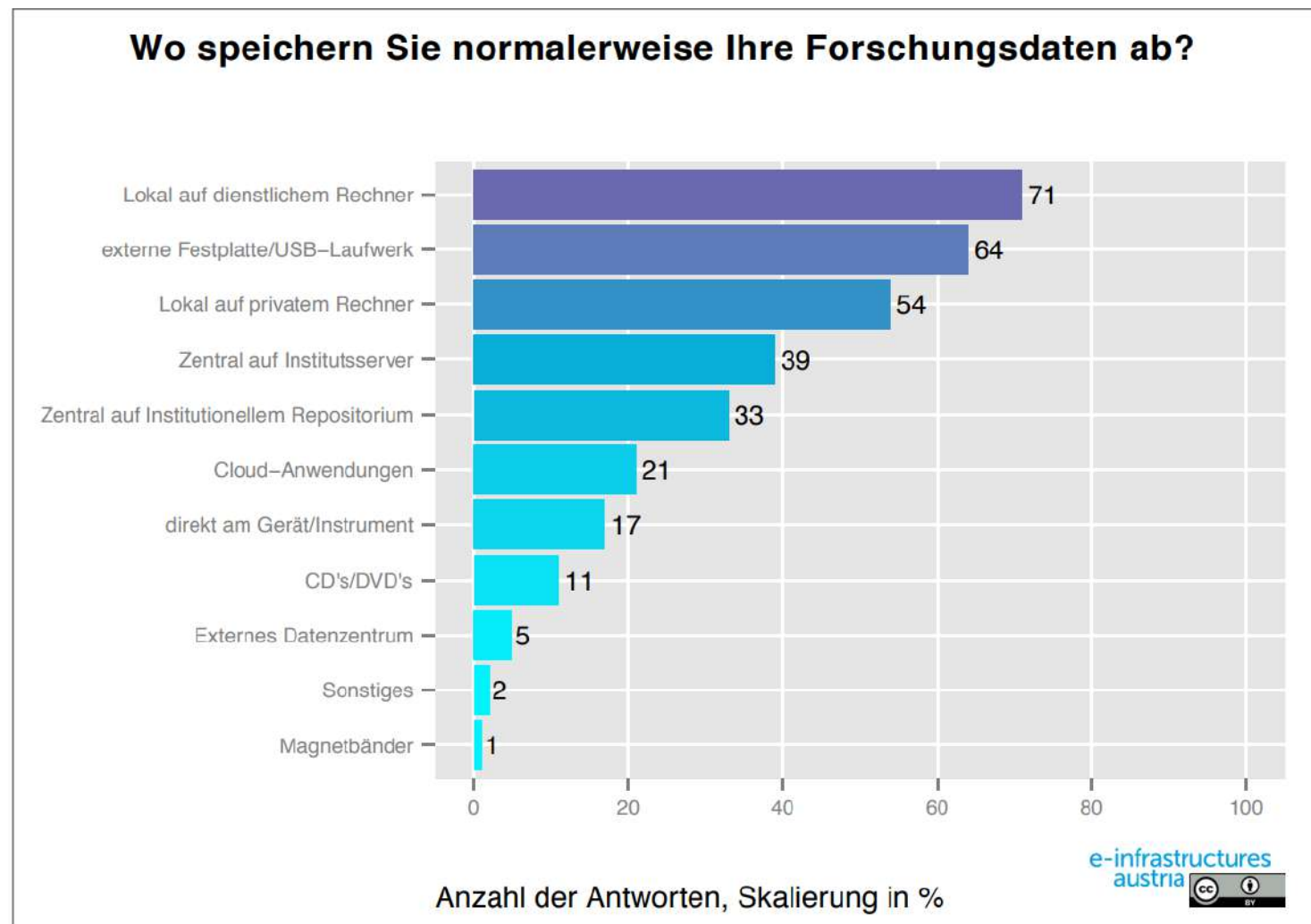
slido

Where do you keep your data?



 Start presenting to display the poll results on this slide.

Managing data during research



Managing data during research

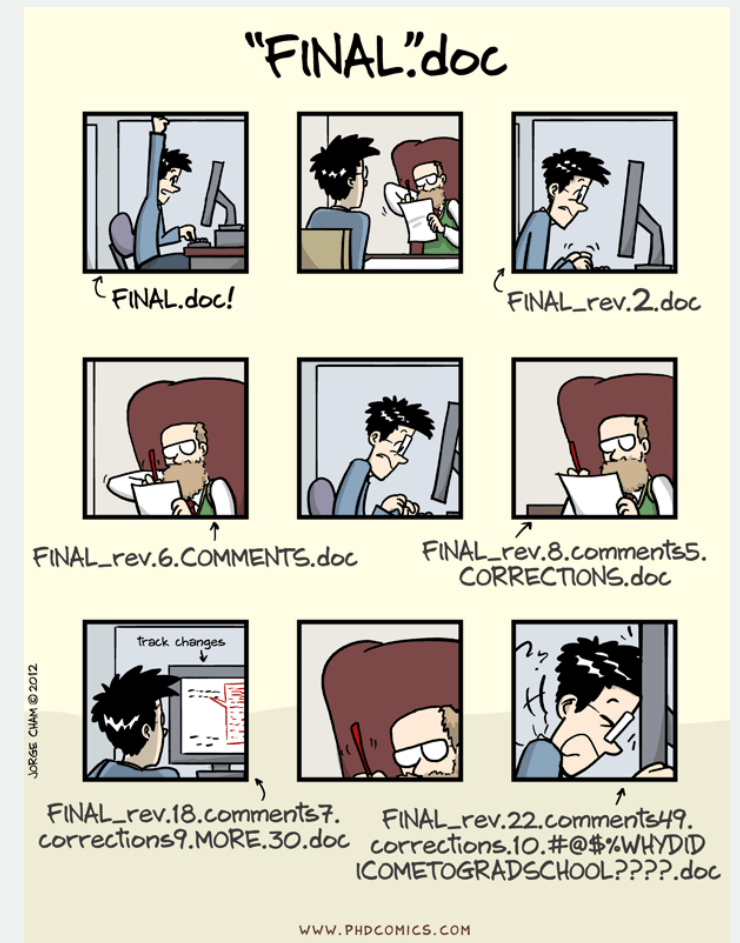
If you loose your data there will be nothing to share!

Recreating or recollecting data can be

- impossible
 - e.g. observational data
- too expensive
 - e.g. cost of computational power

How do you manage data during the project?

- **file naming convention**
- versioning
- backups
- should the access be restricted?
- who is responsible?





Data sharing - example

Code and data are hosted in a public git repository on GitHub.

Read access is open to everyone. Write-access is limited to the researchers working on the project.

Permissions are managed via Github's account system using SSH keys.

Backup vs archiving and preservation (traditional view)



Data managed during the project

- Changed/deleted
- Backup



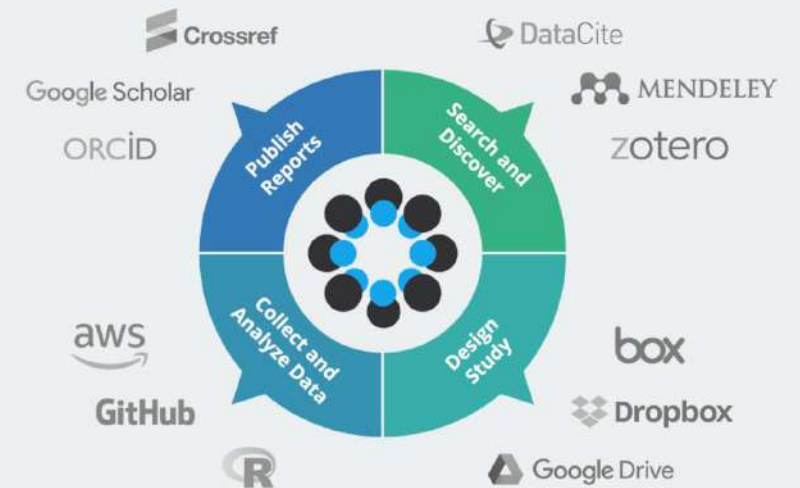
Stable snapshot of data

- Moved into a repository
- Enriched with metadata and licensing details
- Not only backups

What makes a system a repository?

Backup vs archiving and preservation (new approaches)

- No need to differentiate between project and post-project phases
- One system can be used for managing and preserving data



<https://www.cos.io/products/osf>

Archiving and preservation

Which data will be shared?


- What has to be kept?
- What can't be recreated?
- What is potentially useful to others?
- What legally must be destroyed?

Where will the data be deposited?

- not all of the data must be shared in the same way

Are there any embargo periods?

For how long?

What is the cost and who will pay for it? 

Which license to use?

Preservation, repositories,
costs are discussed
in separate lectures!

Archiving and preservation - example



The following files are relevant to reproduce the experiment and should be preserved

- *README.md* – Text file containing instructions on how to run the experiment
- Both *Jupyter notebooks* - The experiment's code and documentation
- *Dockerfile* - To build a docker container for running the experiment
- *requirements.txt* - List of python dependencies required by the experiment
- *documentation/architecture.png* - Architectural diagram of the experiment
- *documentation/description.txt* – Text file describing the correlation plot's content
- *documentation/metadata.xml* - Metadata relevant to the experiment

Note: input data was not selected for preservation

- it is maintained by existing repository (easy to get)

Where to find a repository?

1. Use Domain specific repository

- e.g. chEMBL (if you work with molecules)

2. Use Institutional repository

- e.g. phaidra.univie.ac.at (if you work at Uni Wien)

3. Search registry to find a relevant one

- e.g. re3data.org

4. Use *catch-all* repository

- e.g. zenodo.org

re3data.org


Repository details


TU Data

General Institutions Terms Standards

Name of repository	TU Data
Repository URL	https://researchdata.tuwien.ac.at/
Description	TU Data is an institutional repository of TU Wien to enable storing, sharing and publishing of digital objects, in particular research data. It facilitates the funders' requirements for open access to research data and the FAIR principles by making research output findable, accessible, interoperable and re-usable. This service is developed by the TU Wien Center for Research Data Management and hosted by TU.it.
Content type(s)	Standard office documents Archived data
Keyword(s)	FAIR interdisciplinary
Repository type(s)	institutional
Mission statement for designated community	https://researchdata.tuwien.ac.at/
Research data repository language(s)	English
Data and/or service provider	service provider data provider

[Back to search](#) [Submit a change request](#) [Get a badge](#)


 Cite this re3data.org record:
re3data.org: TU Data; editing status 2021-06-02; re3data.org - Registry of Research Data Repositories.
<http://doi.org/10.17616/R31NJMYD> last accessed: 2022-03-25

re3data.org Search Browse Suggest Resources Contact 

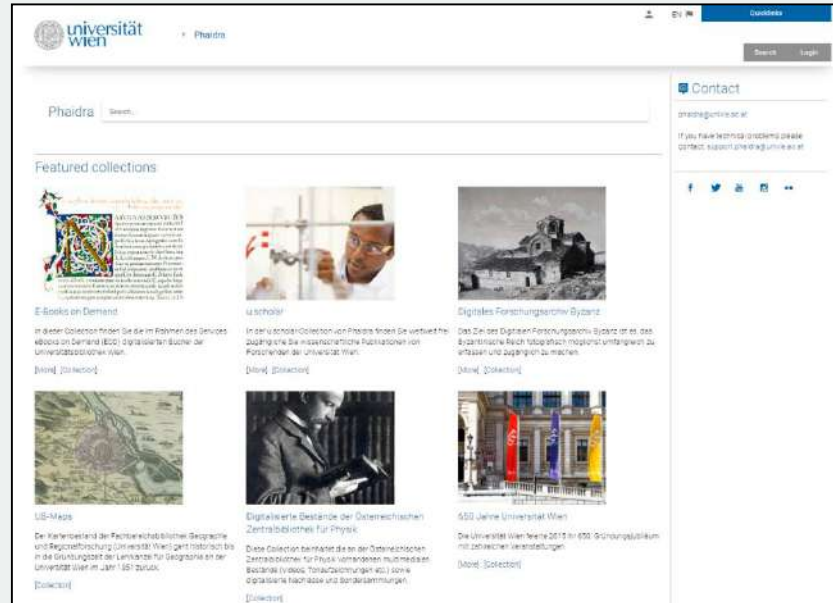
Browse by subject

Graphical Text

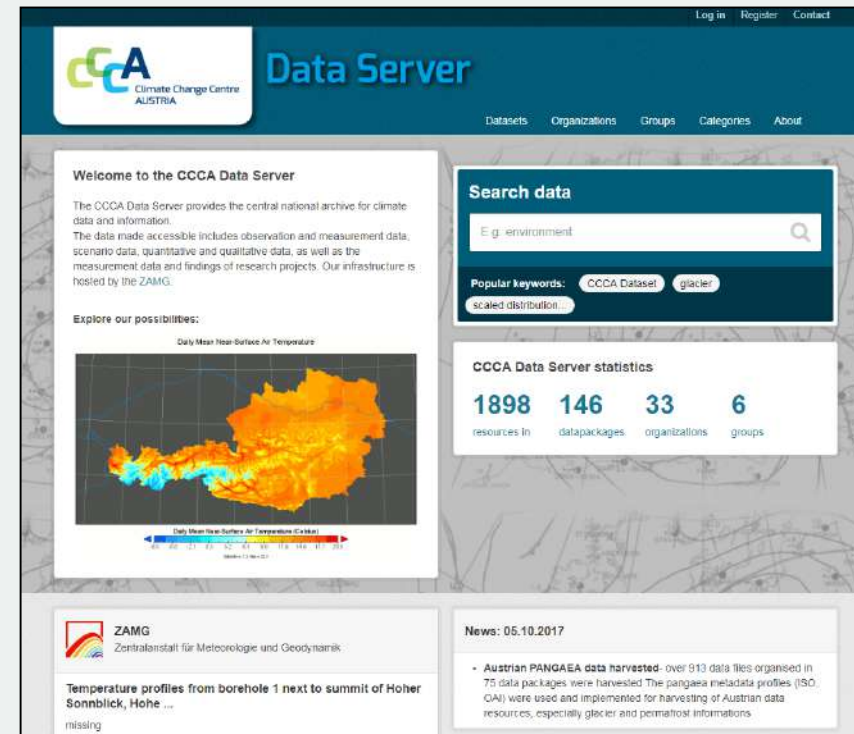
click to zoom into subjects or to select a bottommost subject in the hierarchy as filter for the re3data search page OR click on a top subject to select it as filter



Repositories in Austria - examples



<https://phaidra.univie.ac.at>



<https://data.ccca.ac.at>

V. LEGAL AND ETHICAL ASPECTS

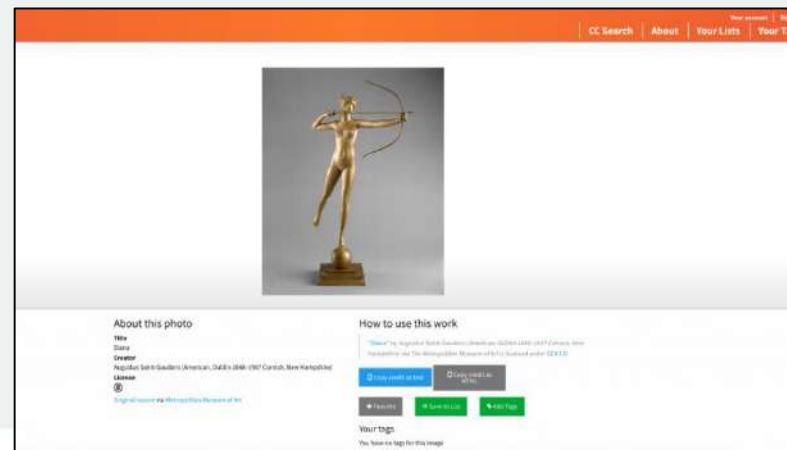
Licenses



CC-0

- waives creator rights -> public domain
- allows anyone to use, re-use, and remix a work without restriction
- Example: all images from Metropolitan Museum of Art in New York

<https://creativecommons.org/2017/02/07/met-announcement/>



Creative Commons

CC-BY (Attribution)

- allows anyone to use, re-use, and remix a work without restriction, also commercially
- You must give appropriate credit, provide a link to the license, and indicate if changes were made.



CC BY-SA (Attribution – ShareAlike)

- all new works must carry the same license



CC BY-ND (Attribution- NoDerivs)

- reuse, but no changes



CC BY-NC

- no commercial use

CC BY-NC-SA

CC BY-NC-ND



Software Licenses

Choose correct license for your software

- Apache, MIT, GNU, BSD, ...

Check licenses of libraries you reuse in your software

- Example: GNU GPL vs GNU LGPL
 - GPL enforces the reusing software to be GPL (also public)
 - LGPL code must be clearly marked, rest of the software can have different license (can be private)

Software licenses can also be used for data

Choose an open source license

Which of the following best describes your situation?

- I want it simple and permissive.**
The MIT License is a permissive license that is short and to the point. It lets people do anything they want with your code as long as they provide attribution back to you and don't hold you liable.
Babel, .NET Core, and Rails use the MIT License.
- I'm concerned about patents.**
The Apache License 2.0 is a permissive license similar to the MIT License, but also provides an express grant of patent rights from contributors to users.
Elasticsearch, Kubernetes, and Swift use the Apache License 2.0.
- I care about sharing improvements.**
The GNU GPL v3 is a copyleft license that requires anyone who distributes your code or a derivative work to make the source available under the same terms, and also provides an express grant of patent rights from contributors to users.
Ansible, Bash, and GIMP use the GNU GPL v3.

What if none of these work for me?

- My project isn't software.**
There are licenses for that.
- I want more choices.**
More licenses are available.
- I don't want to choose a license.**
You don't have to.

The content of this site is licensed under the Creative Commons Attribution 3.0 Unported License. About Terms of Service Contact with +3 by GitHub, Inc. and You!

<https://choosealicense.com>

License💬 example



The external datasets are using permissible licenses which allows us the usage and redistribution of the following data:

- * Ufo Sightings - Creative Commons Attribution 4.0
- * Alcohol Consumption - Free to use and distribute according to <http://www.oecd.org/termsandconditions/> -Section C - Permitted use

All code, data and documentation is available on Github and is licensed under the MIT license.

More on legal aspects of data management

Verena Dolovai will give a lecture in May



SUMMARY

Tips for writing DMPs

DMP can reveal how solid your work is

Seek advice - consult and collaborate

When answering questions from checklists write coherent text

Be specific when referring to tools and standards

Assign responsibilities and name responsible personnel

Tips for writing DMPs

Think about things early...

- Negotiation on licenses and consent agreement may preclude later sharing if not careful
- Manage your data correctly from the very beginning
 - backups, file naming conventions, access restrictions, metadata collection
- Plan your budget

Decisions made early on affect what you can do later

DMPs are not that perfect

Data Management Plans

- are manually created
- depend on scientific honesty
- focus mainly on input and output data
- provide very general overview of the experiment
- have scarce information about the process
- cannot be automatically validated
- do not support sufficiently the reproducibility of research

Lecture on maDMPs will focus on how to fix some of these problems

You should know and be able to explain

- Why and what for we need DMPs
- how to improve your own data management
- what a DMP is and what kind of information it contains
 - Data
 - Metadata
 - Repositories
 - Licensing
 - ...

Useful resources

Center for RDM at TUW (login first!)

- <https://www.tuwien.at/en/research/rti-support/research-data/center-for-rdm>

Managing and sharing data by UK Data Archive

- <http://www.data-archive.ac.uk/media/2894/managingsharing.pdf>

EUDAT webinars on data management

- <https://www.eudat.eu/events/webinar/research-data-management-an-introductory-webinar-from-openaire-and-eudat>

FFG-Akademie: Der Data Management Plan (DMP) in Horizon 2020 (Webinar)

- https://www.ffg.at/europa/veranstaltungen/ffg-akademie_2017-10-18

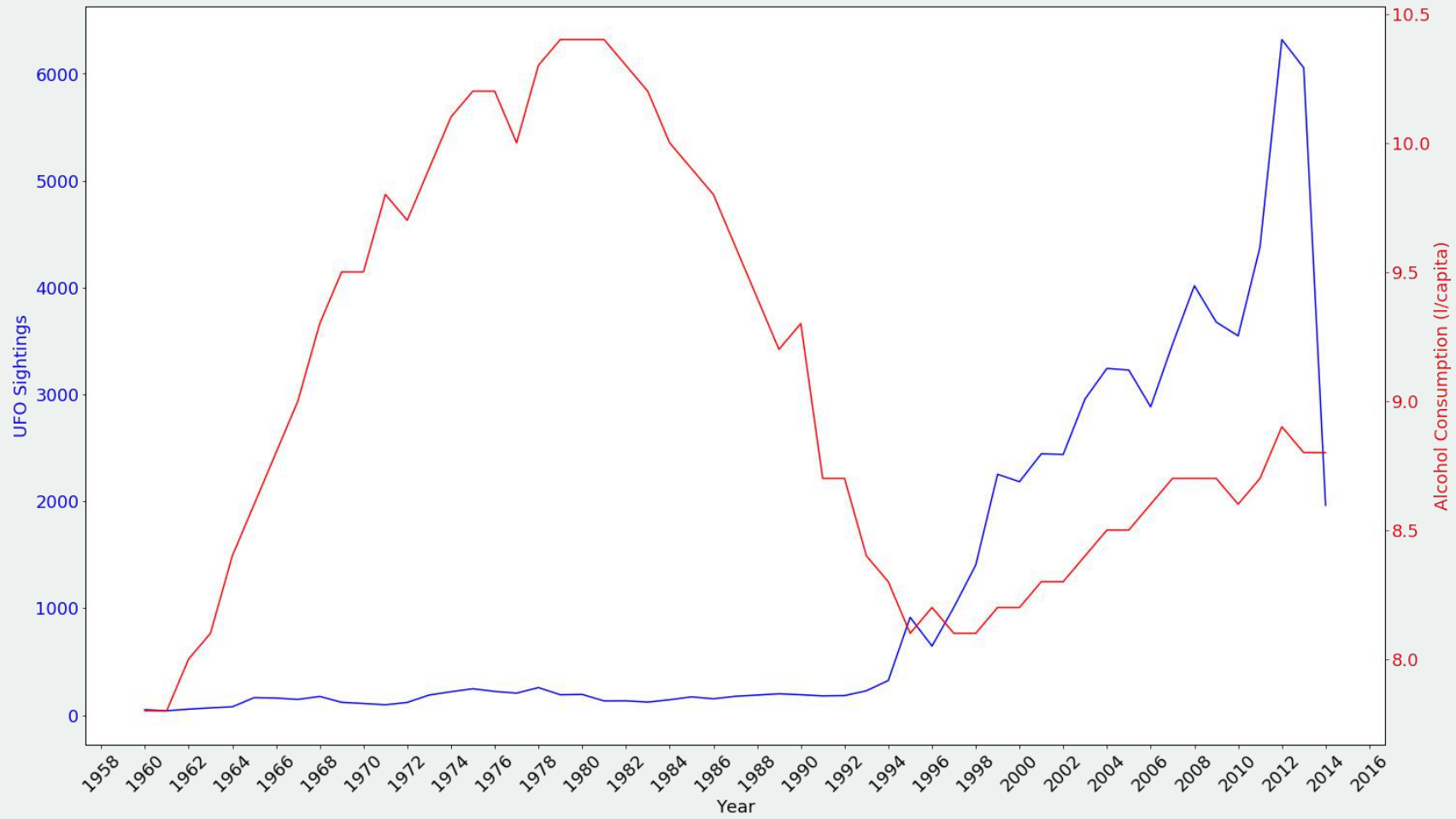
Ten Simple Rules

- <http://dx.doi.org/10.1371/journal.pcbi.1004525>

DMP Checklist

- http://www.dcc.ac.uk/sites/default/files/documents/resource/DMP/DMP_Checklist_2013.pdf

Correlating Alcohol Consumption and UFO Sightings in the USA



<https://github.com/mdietrichstein/digitalpreservation-dmp>



Enabling Precise Identification and Citability of Dynamic Data

Andreas Rauber

Vienna University of Technology
Favoritenstr. 9-11/188
1040 Vienna, Austria
rauber@ifs.tuwien.ac.at
<http://ww.ifs.tuwien.ac.at/~andi>



FACULTY OF **INFORMATICS**

Outline

-
- Why should we want to cite data?
 - What identifier system should I use?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - Summary
-

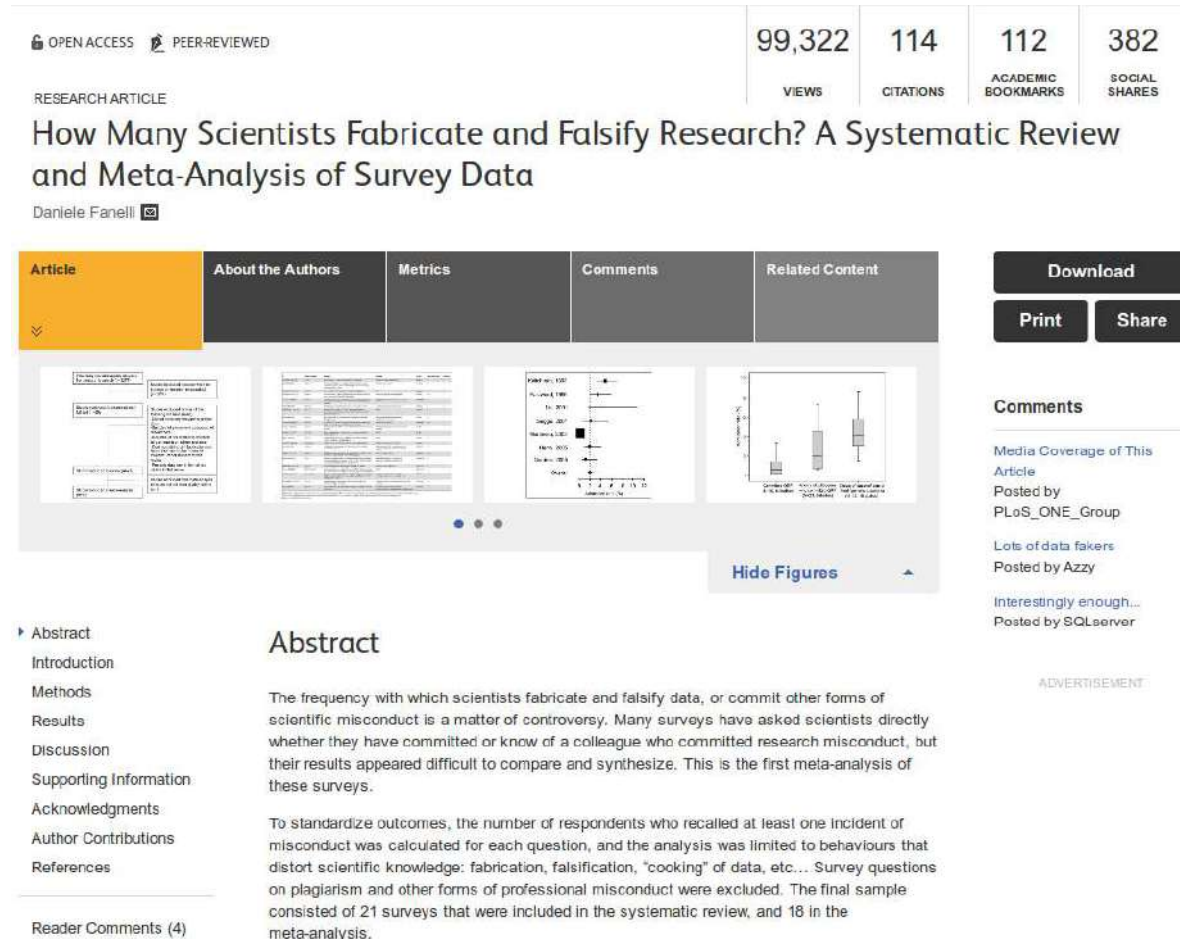
Why to cite data?

- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...

- Why should we cite data?
 - Prevent scientific misconduct (“extrinsic”) ?

Prevent Scientific Misconduct

- On average: ~ 2% of scientists admitted they had "fabricated" (made up), "falsified" or "altered" data to "improve the outcome" at least once.
- 34% admitted to other questionable research practices including "failing to present data that contradict one's own previous research" and "dropping observations or data points from analyses based on a gut feeling that they were inaccurate."
- 14% knew someone who had fabricated, falsified or altered data.
- Up to 72% knew someone who had committed other questionable research practices.



OPEN ACCESS PEER-REVIEWED
 99,322 VIEWS 114 CITATIONS 112 ACADEMIC BOOKMARKS 382 SOCIAL SHARES
 RESEARCH ARTICLE
How Many Scientists Fabricate and Falsify Research? A Systematic Review and Meta-Analysis of Survey Data
 Daniele Fanelli

Article About the Authors Metrics Comments Related Content
 Download Print Share
 Comments
 Media Coverage of This Article
 Posted by PLoS_ONE_Group
 Lots of data fakers
 Posted by Azzy
 interestingly enough...
 Posted by SQLserver
 ADVERTISEMENT

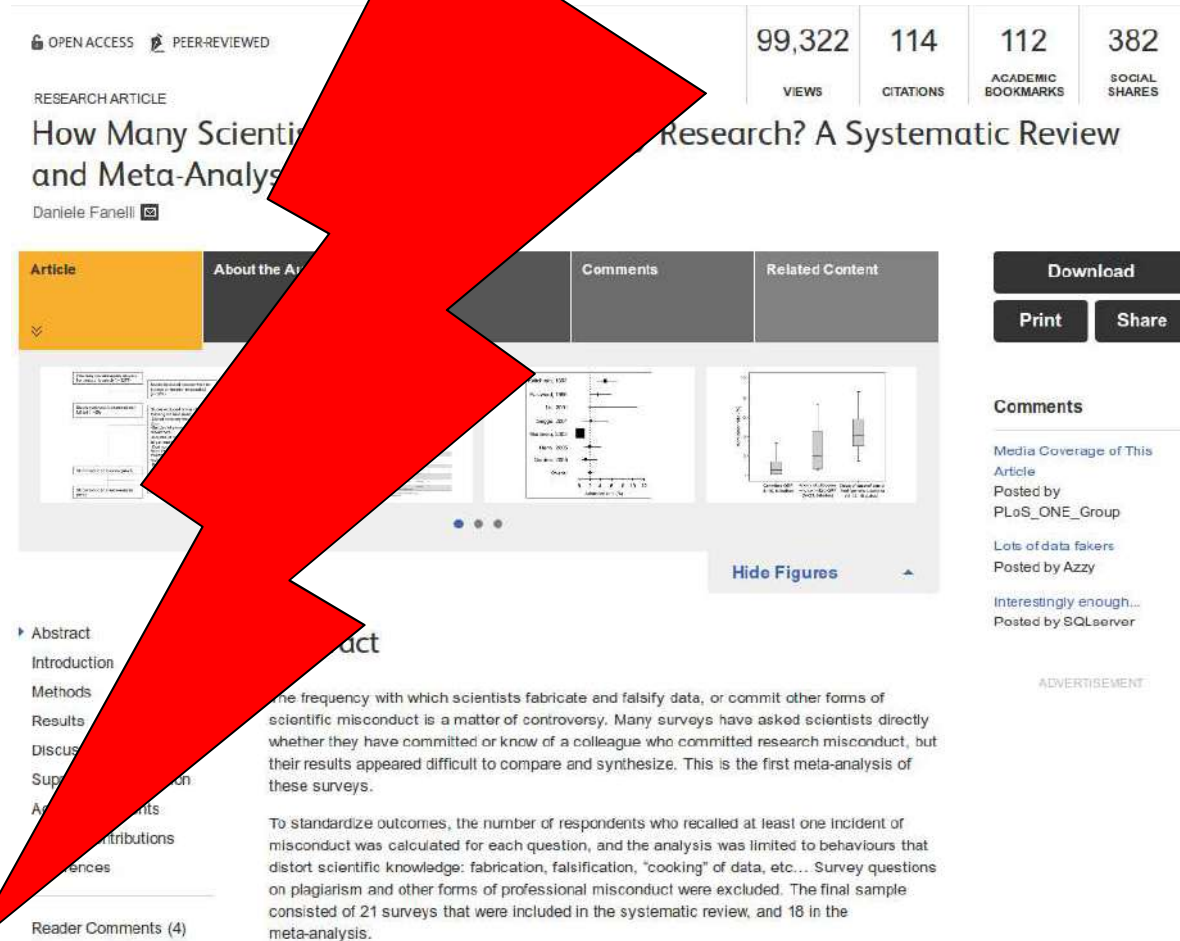
Abstract Introduction Methods Results Discussion Supporting Information Acknowledgments Author Contributions References
 Reader Comments (4)

Abstract
 The frequency with which scientists fabricate and falsify data, or commit other forms of scientific misconduct is a matter of controversy. Many surveys have asked scientists directly whether they have committed or know of a colleague who committed research misconduct, but their results appeared difficult to compare and synthesize. This is the first meta-analysis of these surveys.
 To standardize outcomes, the number of respondents who recalled at least one incident of misconduct was calculated for each question, and the analysis was limited to behaviours that distort scientific knowledge: fabrication, falsification, "cooking" of data, etc... Survey questions on plagiarism and other forms of professional misconduct were excluded. The final sample consisted of 21 surveys that were included in the systematic review, and 18 in the meta-analysis.

Source: <http://www.plosone.org>

Prevent Scientific Misconduct


- On average: ~ 2% of scientists admitted they had "fabricated" (made up), "falsified" or "altered" data to "improve the outcome" at least once.
- 34% admitted to other questionable research practices including "failing to present data that contradict one's own previous research" and "dropping observations or data points from analyses based on a gut feeling that they were inaccurate."
- 14% knew someone who had fabricated, falsified or altered data.
- Up to 72% knew someone who had committed other questionable research practices.



Source: <http://www.plosone.org>

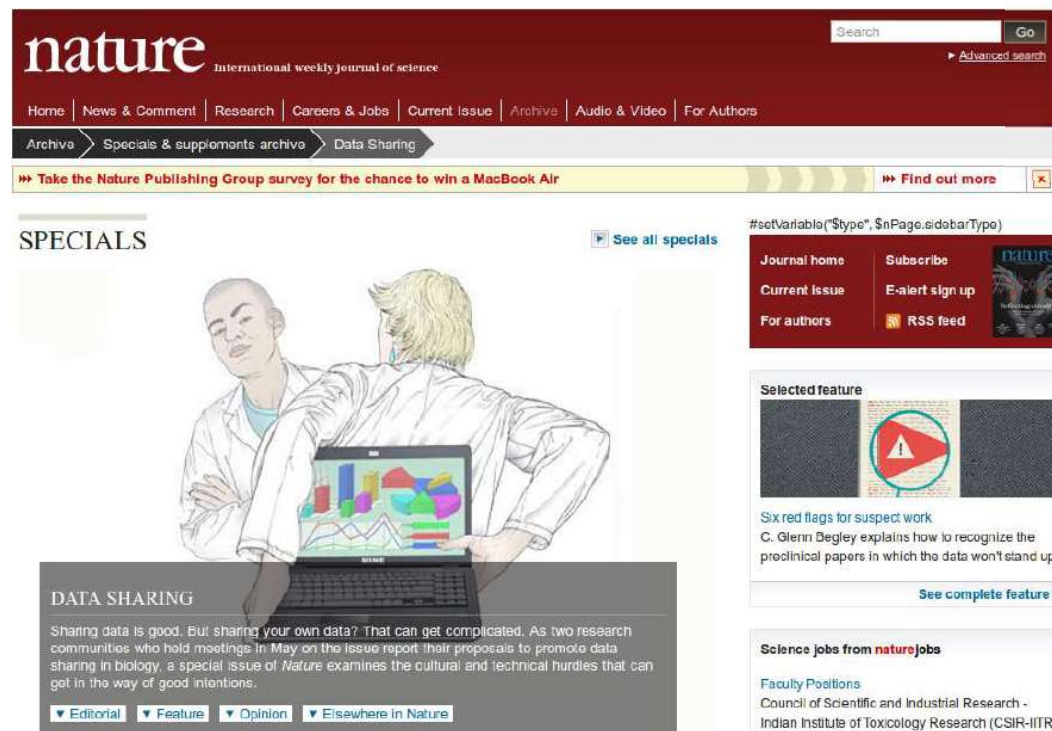
Why to cite data?

- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...

- Why should we cite data?
 - Prevent scientific misconduct (“extrinsic”) ? 
 - Give credit (“altruistic”) ?

Giving credit

- Prime motivator for sharing data
- Shared data gets cited more frequently
- Citations are the currency of science



The screenshot shows the Nature journal website. The main navigation bar includes 'Home', 'News & Comment', 'Research', 'Careers & Jobs', 'Current Issue', 'Archive', 'Audio & Video', and 'For Authors'. Below this, there are links for 'Archive', 'Specials & supplements archive', and 'Data Sharing'. A banner for a survey is visible: 'Take the Nature Publishing Group survey for the chance to win a MacBook Air'. The 'SPECIALS' section features an illustration of two scientists looking at a laptop displaying a colorful bar chart. Below the illustration is a 'DATA SHARING' article with the text: 'Sharing data is good. But sharing your own data? That can get complicated. As two research communities who held meetings in May on the issue report their proposals to promote data sharing in biology, a special issue of Nature examines the cultural and technical hurdles that can get in the way of good intentions.' The article has tags for 'Editorial', 'Feature', 'Opinion', and 'Elsewhere in Nature'. On the right side, there is a sidebar with links for 'Journal home', 'Subscribe', 'Current Issue', 'E-alert sign up', 'For authors', and 'RSS feed'. Below this is a 'Selected feature' section with a play button icon and the text: 'Six red flags for suspect work. C. Glenn Degley explains how to recognize the preclinical papers in which the data won't stand up.' At the bottom of the sidebar, there is a 'Science jobs from naturejobs' section with links for 'Faculty Positions' and 'Council of Scientific and Industrial Research - Indian Institute of Toxicology Research (CSIR-IITR)'.

Giving credit

- Prime motivator for sharing data
- Shared data gets cited more
- Citations are the currency of science



Why to cite data?

- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...

- Why should we cite data?
 - Prevent Scientific misconduct (“extrinsic”) ?
 - Give credit (“altruistic”) ?
 - Show solid basis (“egoistic”) ?



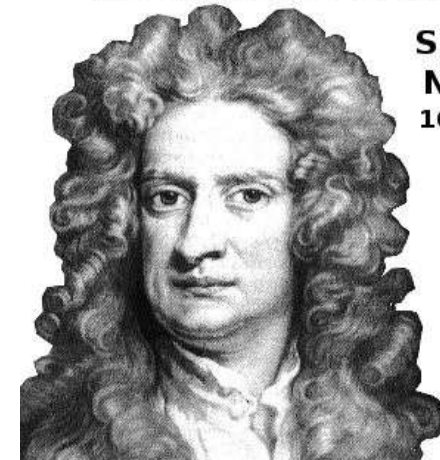
Citing to give credit

Why do we cite papers? (“related work”)

- Fundamental basis for own work – foundation!
- No need to prove - it’s been done!
- Speed-up the process, efficiency
- Basis for discourse, scientific work, ...






"If I have seen further, it has been by standing on the shoulders of giants."



**Sir Isaac
Newton**
1643-1727

Why to cite data?

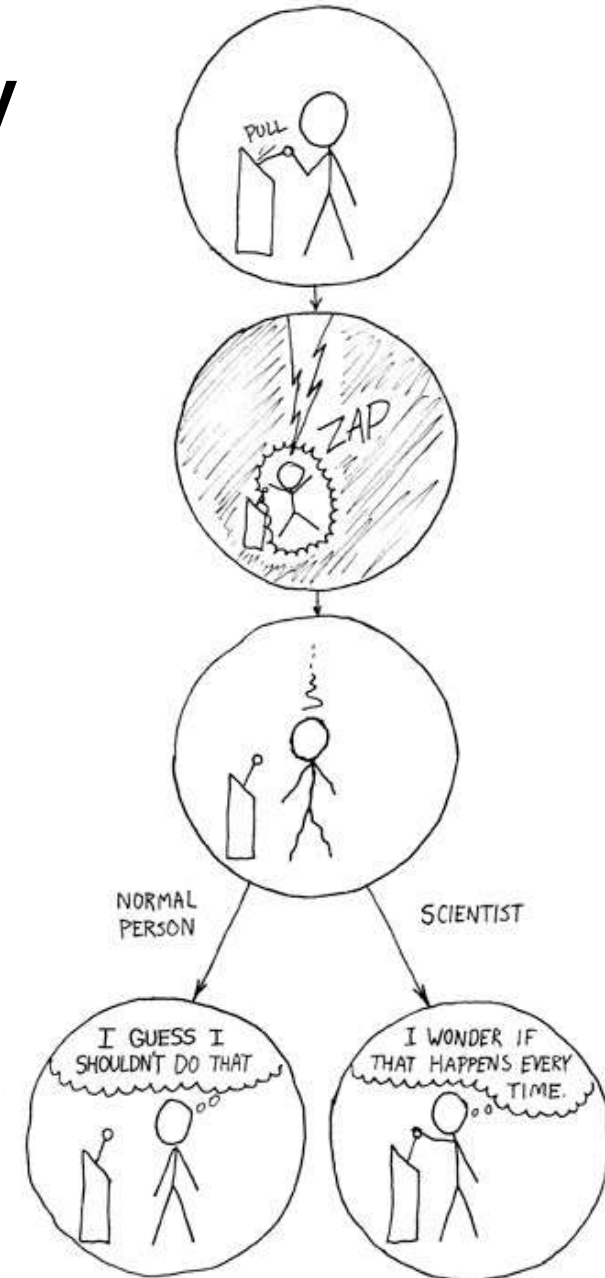
- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...

- Why should we cite data?
 - Prevent Scientific misconduct (“extrinsic”) ? 
 - Give credit (“altruistic”) ? 
 - Show solid basis (“egoistic”) ? 
 - Enable **reproducibility**, re-use (extrinsic + altruistic + egoistic) ?



Reproducibility

- Reproducibility is core to the scientific method
- Focus not on misconduct – but on complexity and the will to produce good work
- Should be easy
 - Get the code, compile, run, ...
 - Why is it difficult?




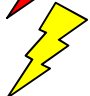


Reproducibility

- requires
 - Transparency, requires
 - Documentation, provides
 - Context, requires
 - » **Citation**, requires
 - » **Identification**
- Increases impact
- Increases trust
- Fosters reuse






<https://xkcd.com/978/>

Why to cite data?

- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...

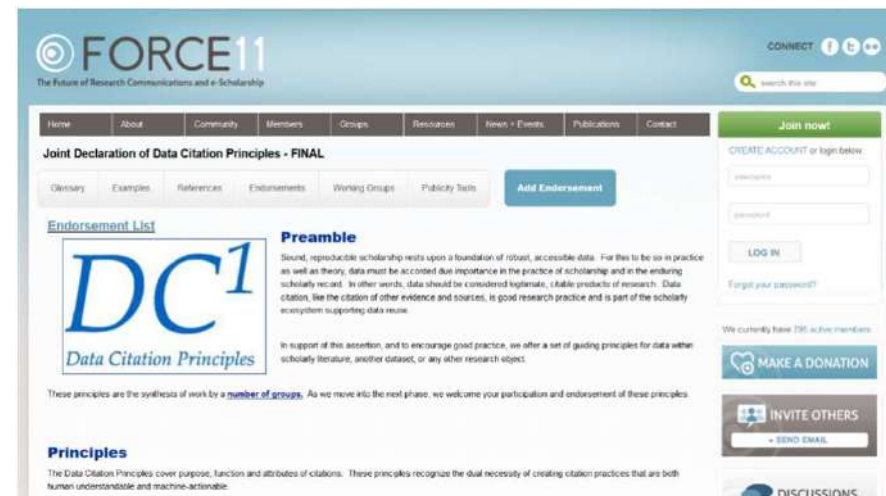
- Why should we cite data?
 - Prevent Scientific misconduct (“extrinsic”) ? 
 - Give credit (“altruistic”) ? 
 - Show solid basis (“egoistic”) ? 
 - Enable reproducibility, re-use (extrinsic + altruistic + egoistic) 

Why to cite data?

- Data is the basis for almost everything
 - eScience, digital humanities,
 - Industry 4.0
 - Driving policies, society, ...
- Why should we cite data?
 - Prevent Scientific misconduct (“extrinsic”) ? 
 - Give credit (“altruistic”) ? 
 - Show solid basis (“egoistic”) ? 
 - Enable reproducibility, re-use (extrinsic + altruistic + egoistic) ? 
 - **Because it’s what you do if you do good work, speeding up the process of scientific discovery, efficiency! (“intrinsic”)** 

Joint Declaration of Data Citation Principles

- 8 Principles created by the Data Citation Synthesis Group
- <https://www.force11.org/datacitation>
- The Data Citation Principles cover purpose, function and attributes of citations
- Goal: Encourage communities to develop practices and tools that embody uniform data citation principles



1) Importance

Data should be considered legitimate, citable products of research. Data citations should be accorded the same importance as publications. 

2) Credit and Attribution

Data citations should facilitate giving credit and normative and legal attribution to all contributors to the data.

3) Evidence

Whenever and wherever a claim relies upon data, the corresponding data should be cited.


4) Unique Identification

A data citation should include a persistent method for identification that is machine actionable, globally unique, and widely used by a community.


5) Access

Data citations should facilitate access to the data themselves and to such associated metadata, documentation, code, and other materials, as are necessary for both humans and machines to make informed use of the referenced data.

6) Persistence

Unique identifiers, and metadata describing the data, and its disposition, should persist - even beyond the lifespan of the data they describe. 

7) Specificity and Verifiability

Data citations should facilitate identification of, access to, and verification of the specific data that support a claim. Citations or citation metadata should include information about provenance and fixity sufficient to facilitate verifying that the specific time slice, version and/or granular portion of data retrieved subsequently is the same as was originally cited. 

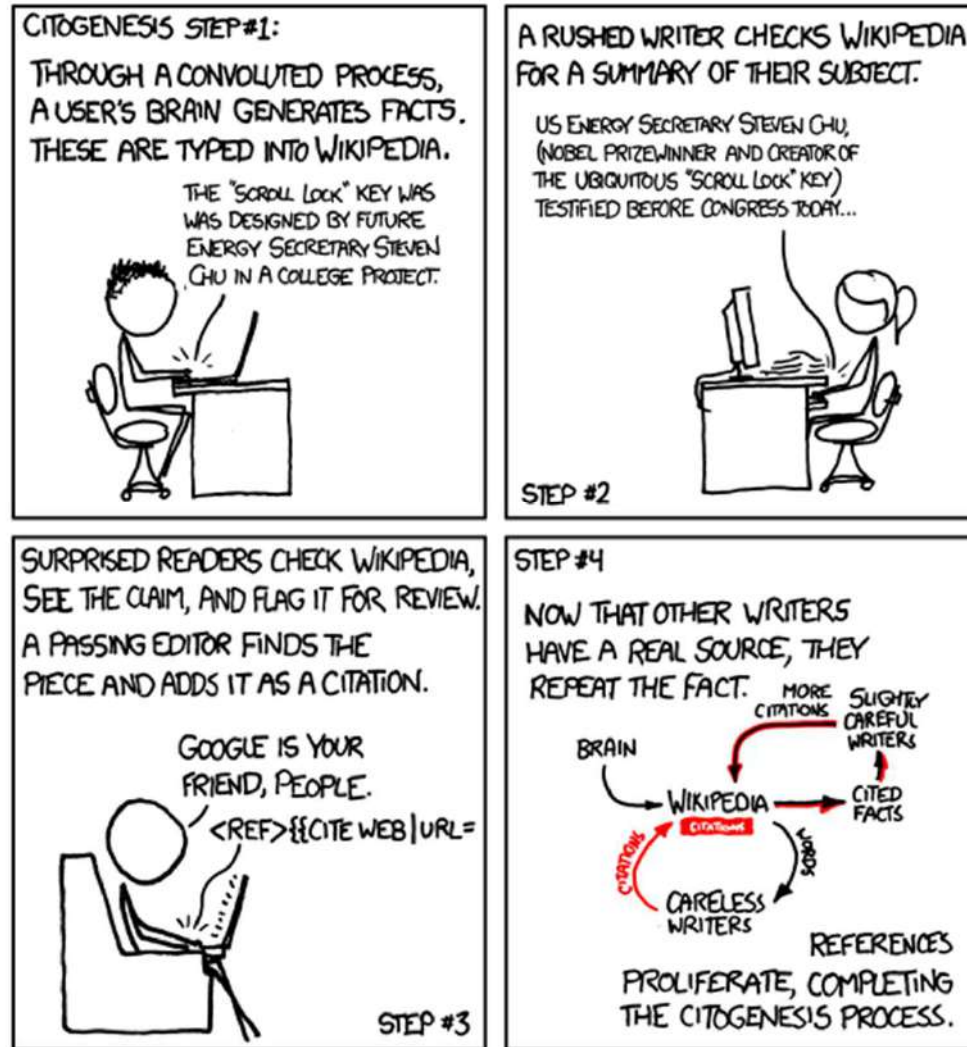
8) Interoperability and flexibility

Data citation methods should be sufficiently flexible to accommodate the variant practices among communities, but should not differ so much that they compromise interoperability of data citation practices across communities.

Benefits of Citation

- Identification
- Documentation
- Context
- Impact
- Transparency
- Reproducibility
- Reuse

WHERE CITATIONS COME FROM:



<https://xkcd.com/978/>



Standard Elements of Data Citation

- Classical bibliographic details:
 - Author, date, edition
 - Publisher, version
- Specific details:
 - Feature name, resource type
 - Unique numeric fingerprint (hash)
 - Persistent identifier
 - Location
- But there is more to it...
Landing pages – “end credits in movies”

Outline

-
- Why should we want to cite data?
 - What identifier system should I use?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - Summary
-

Identifiers

- Identifier is a symbol that uniquely identifies an object.
 - Used to identify (digital) objects
 - References the location
 - Provides metadata
 - Can be resolved
 - Several identifier types exist



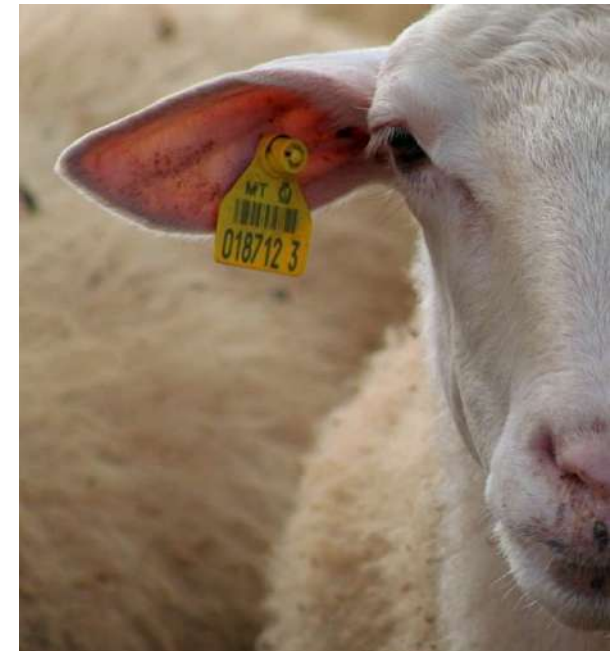
Traditional Mechanisms

- International Standard Serial Number (ISSN)
 - Unique eight-digit number
 - Identifiers periodical publications
 - Can be encoded as URN
- International Standard Book Number (ISBN)
 - Unique commercial book identifier barcode
 - 13 (since 2007) or 10 digits with checksum
 - ISBN-10: 3836217155
 - ISBN-13: 978-3836217156



Unique Identifiers for Digital Objects

- Originally:
 - Uniform Resource Name (URN)
 - Uniform Resource Locator (URL)
 - Uniform Resource Characteristic (URC, metadata, replaced by RDF)
- Uniform Resource Identifier (URI)
 - Encompasses URN and URL
 - Can be resolvable, but need not be
 - Includes ISBN, etc.
- Delegating Methods
 - Handle System (basis for all others)
 - Digital Object Identifier (DOI)
 - Persistent URL (PURL)
 - Archival Resource Key (ARK)



URLs and Persistency?

- Standard URLs are not forever
 - Describe network locations
 - Not suitable for the long term
 - Link rot:
“half of the links in publications are not available after 5 to 7 years” (precise numbers vary...)

- Solution: persistent identifiers (PIDs)

PURL

- Persistent uniform resource locator
- Developed by Online Computer Library Center in 1995
- Based on HTTP forwarding
 - Only resolution
 - No metadata
- Provides curation and URL resolvers
- Can be hosted on own servers or centrally
- Is free
- Example: <http://purl.fdlp.gov/GPO/gpo49354>
 - [Catalog of U.S. Government Publications](#)

PURL Domains

Logged in as [sproell@sba-research.org](#) ([log out](#))

PURL Domain Administration

[Home](#) [PURLs](#) [Users](#) [Groups](#) **[Domains](#)** [Admin](#) [Help](#)

1) Choose an action to take on domains

Domain administration options.

Create a new domain



2) Create a new domain

Fill in the following information to create a new domain.

Name:

Domain ID:


Maintainer IDs (one per line):

Writer IDs (one per line):

Public? (Applies solely to top-level domains):

Create Successful

status: Pending approval

id: /APC2014 

name: Advanced Practitioner Course 2014

public: false

maintainers: sproell@sba-research.org

writers: sproell@sba-research.org

PURL Registration

Logged in as [sproell@sba-research.org](#) ([log out](#))


PURL Administration

Home **PURLs** Users Groups Domains Admin Help

1) Choose an action to take on PURLs

PURL administration options.

Create a new PURL



2) Create a new PURL

Fill in the following information to create a new PURL.

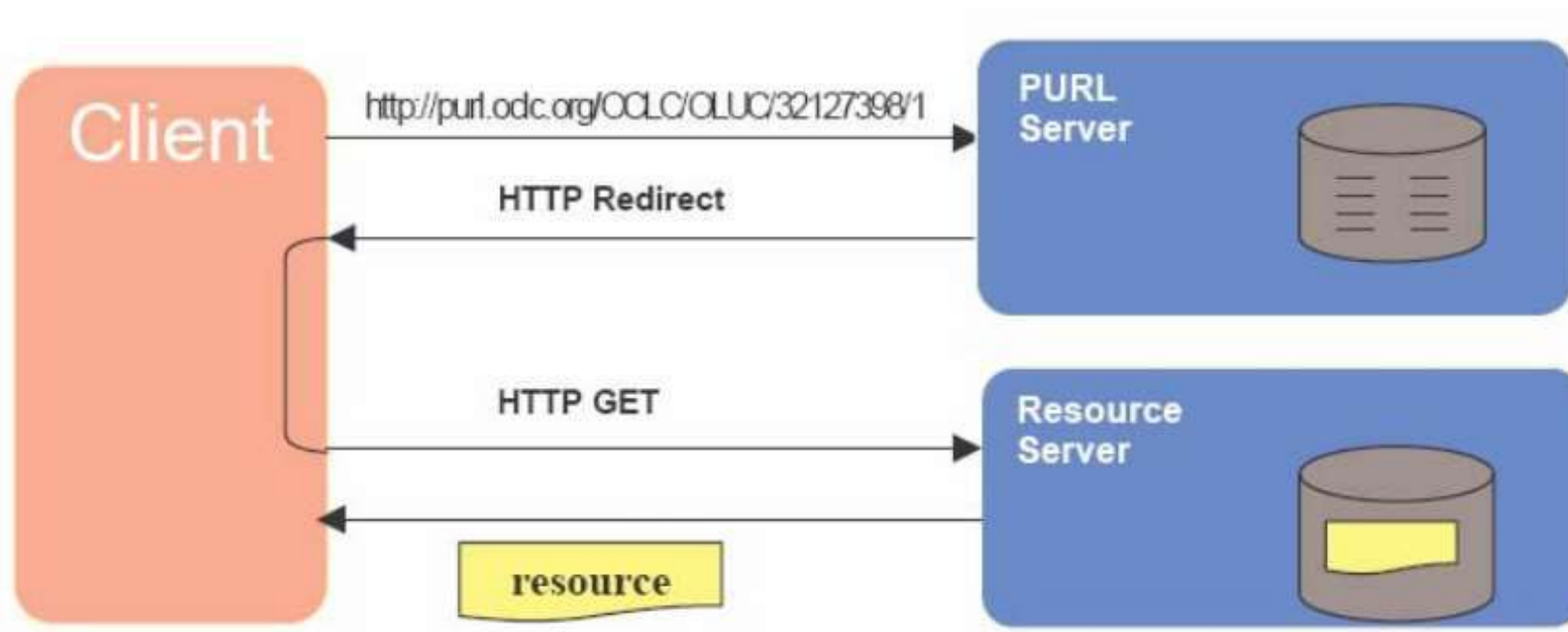
Path:

Target URL:

Maintainers IDs
(one per line):

[Advanced](#)

PURL - Resolution



© MPDL

<http://www.mpi.nl/DAM-LR/meeting5/Persistent%20Identifiers.pdf>

Uniform Resource Identifiers

- Uniform Resource Identifier
 - Name
 - Location (e.g. Web server)
- Combination of
 - namespace identifier (NID) and a
 - namespace specific string (NSS)
- Naming scheme for URNs:
 - urn: <NID> :<NSS>
 - urn:isbn:0451450523
 - urn:ietf:rfc:2648
- Note: Cool URIs don't change!
<https://www.w3.org/Provider/Style/URI>



URN

- Main characteristics and functions of a URN usually include
 - Global scope of names
 - Global uniqueness
 - Persistence
 - Scalability
 - Legacy support
 - Extensibility
 - Independence
 - Resolution (handle system)



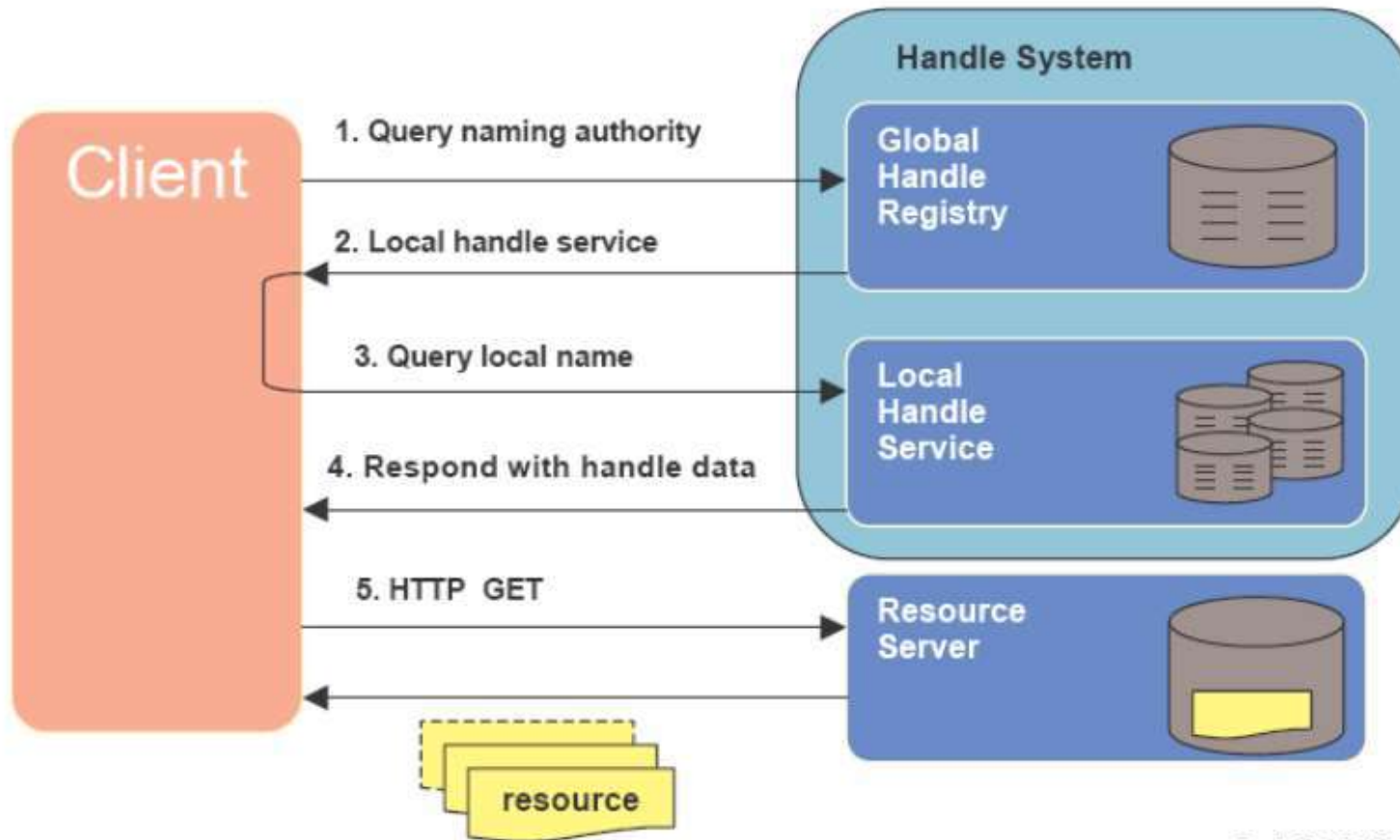
Handle

- Distributed persistent naming system
- Conforms to URN framework
- Used by most identifier systems
- Persistent identifier consists of two parts:
 - Naming authority
 - Name (must be unique string to the authority)
- Digital objects on the Internet can be assigned, managed and resolved by handles
- Resolved by global handle service
- E.g. <http://hdl.handle.net>

Handle

- Main points
 - Handles are unique and persistent
 - Operations on handle system have to be authorized
- Syntax:
 <Handle Naming Authority> ,/‘ <Handle Local Name>
- Example:
 - 10.1045/january2013-burns
- Available Services:
 - <http://hdl.handle.net>



Handle Resolution



© MPDL

[8] <http://www.mpi.nl/DAM-LR/meeting5/Persistent%20Identifiers.pdf>

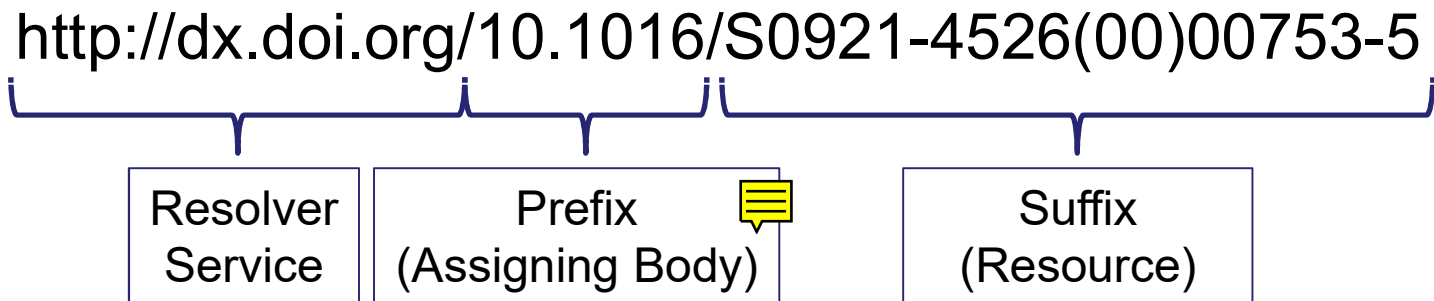
Digital Object Identifier (DOI)

- **Digital Identifier of an Object**
 - not "Identifier of a Digital Object" 
 - "click on it and do something"
- Identifier scheme administered by the International DOI Foundation (IDF) 
- Relies on the handle concept
- Provides an actionable, interoperable, persistent link
- International Standard: ISO 26324 (May 2012)
- TU Wien Library provides DOI Services for Austria
https://www.ub.tuwien.ac.at/pid/pid_result.html




Digital Object Identifier (DOI)

- Consists of three parts:



- Resource can be any entity (thing: physical, digital, or abstract)
- DOI: 10.1594/PANGAEA.724325
- Resolver services lead to landing page
 - <http://dx.doi.org/>
 - <http://dx.doi.org/10.1594/PANGAEA.724325>

DOI: Guidelines

- Suffix must be unique within the prefix
- Suffix is case insensitive
- UTF-8
- Recommendations:
 - Use short suffixes, people have to type them
 - Do not use special characters if possible
 - Avoid semantics in the suffix string, as its semantics could change (“**no semantics in an identifier!**”) 
 - Slightly contradicted in “fragment identifiers”
(*personal comment: avoid!*)

Metadata:

- DOI Kernel Metadata

https://www.doi.org/doi_handbook/4_Data_Model.html

- Other identifiers (isbn, issn, ...)
- structural types (e.g. *physical, digital, performance, abstraction*)
- modes (*audio, visual, tangible, olfactory, tasteable, none*),
- linkedCreation,
- linkedParty, date of birth/death, territory, ...
- (several more)

- DOI Data Dictionary

https://www.doi.org/doi_handbook/schemas/dd/intro.html


Example: Formatted Citations:

- `curl -LH "Accept: text/x-bibliography; style=apa"`
<http://dx.doi.org/10.1126/science.169.3946.635>
 - Frank, H. S. (1970). The Structure of Ordinary Water: New data and interpretations are yielding new insights into this fascinating substance. *Science*, 169(3946), 635–641.
doi:10.1126/science.169.3946.635
- `curl -LH "Accept: text/x-bibliography; style=bibtex"`
<http://dx.doi.org/10.1126/science.169.3946.635>
 - `@article{Frank_1970, title={The Structure of Ordinary Water: New data and interpretations are yielding new insights into this fascinating substance}, volume={169}, ISSN={1095-9203}, url={http://dx.doi.org/10.1126/science.169.3946.635}, DOI={10.1126/science.169.3946.635}, number={3946}, journal={Science}, publisher={American Association for the Advancement of Science (AAAS)}, author={Frank, H. S.}, year={1970}, month={Aug}, pages={635–641}}`

DataCite Metadata Store

Metadata Store [Search](#) [Schema](#) [OAI-PMH](#) [Content Resolver](#) [Stats](#) [Handle Server](#)

This service is for testing only.













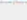
Dataset

- [Register new Dataset](#)
- [List all Datasets](#)
- [Find by DOI](#)




View

- [API documentation](#)

▼ List all Datasets

DOI	Is Active	Is Ref Quality	Updated	Minted	Latest Metadata Version	
10.5072/A1WDE5GFFRECBN879	true	false	2014-09-04 14:42 UTC	2014-09-04 14:40 UTC	0 (2014-09-04 14:40:27.0)	
10.5072/ECFA2TZUNBV4562	true	false	2014-09-04 14:19 UTC	2014-09-04 14:16 UTC	0 (2014-09-04 14:16:24.0)	
10.5072/6P76C3PB12345	true	false	2014-09-04 13:57 UTC	2014-09-04 13:57 UTC	0 (2014-09-04 13:57:00.0)	
10.5072/726855ASWWSSFDBNDS	true	false	2014-07-10 08:45 UTC	2014-07-10 08:35 UTC	0 (2014-07-10 08:35:17.0)	
10.5072/FZJK4ZJJNMDN353	true	false	2014-07-09 15:03 UTC	2014-07-09 15:03 UTC	0 (2014-07-09 15:03:34.0)	
10.5072/KJHGFDSA6543	true	false	2014-07-09 10:26 UTC	2014-07-09 10:26 UTC	0 (2014-07-09 10:26:00.0)	
10.5072/TPDL2013TUTORIAL	true	false	2013-09-22 10:02 UTC	2013-09-22 10:01 UTC	0 (2013-09-22 10:01:13.0)	
10.5072/DATASET-TPDL-TEST	true	false	2013-09-18 08:10 UTC	2013-09-18 08:08 UTC	0 (2013-09-18 08:08:39.0)	
10.5072/DATASET-TPDL	true	false	2013-09-17 12:58 UTC	2013-09-17 12:55 UTC	0 (2013-09-17 12:55:30.0)	
10.5072/PROELLA1B2C3D4	true	false	2013-08-30 19:04 UTC	2013-08-30 19:04 UTC	0 (2013-08-30 19:04:27.0)	
10.5072/PROELLA1B2C3	true	false	2013-06-28 12:24 UTC	2013-06-28 12:24 UTC	0 (2013-06-28 12:24:04.0)	

List results per page: 30 [50](#) [100](#) | Page 1 of 1

[Home](#) | Language:    | [Logout](#)

How to Get a DOI

1. Request an account at a DOI registration agency
2. Pay a fee
3. Receive login data and your prefix
4. Establish (“mint”) a DOI suffix to be linked to your object providing the required metadata
5. Start citing

DOI – Facts

- Launched in 2000
- Over 5,000 naming authorities (assigners)
- Over 20,000 DOI name prefixes
- Over 148 million DOI names assigned
 - Grows 16 % per year!
- Over 5 billion DOI resolutions per year
- International Standard: ISO 26324 (May 2012)

DOI Registration Agencies

- Are members of the IDF and entitled to assign and maintain DOIs.
- Examples:
 - DataCite
 - CrossRef
 - Bowker
 - CAL
 - Nielsen BookData
 - TIB
 - OPOCE
 - TU Wien





Helping you to find,
access, and reuse data

- Registration Agency for DOIs
- Non-profit membership organization established 2009
- Aims:
 - Establish easier access to research
 - Increase acceptance of research data
 - Support data archiving that will permit results to be verified and re-purposed for future study.

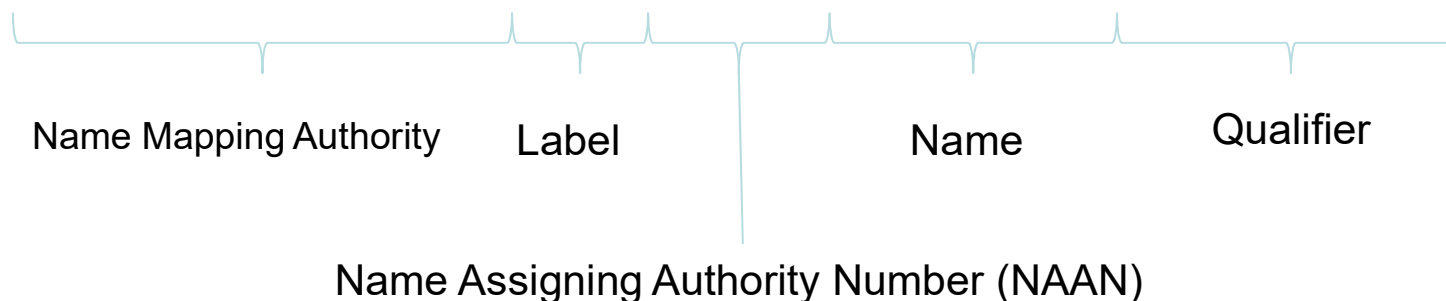
DOI vs. Handle

- Handle only provides the resolution service
- DOI uses the Handle System and adds:
 - **Metadata** (remains even if object is no longer available)
 - Consistency of citations
 - Semantic interoperability (data model)
 - Identification of intellectual property entities
- Used by aggregators, impact factor calculation, ...

Archival Resource Key (ARK)

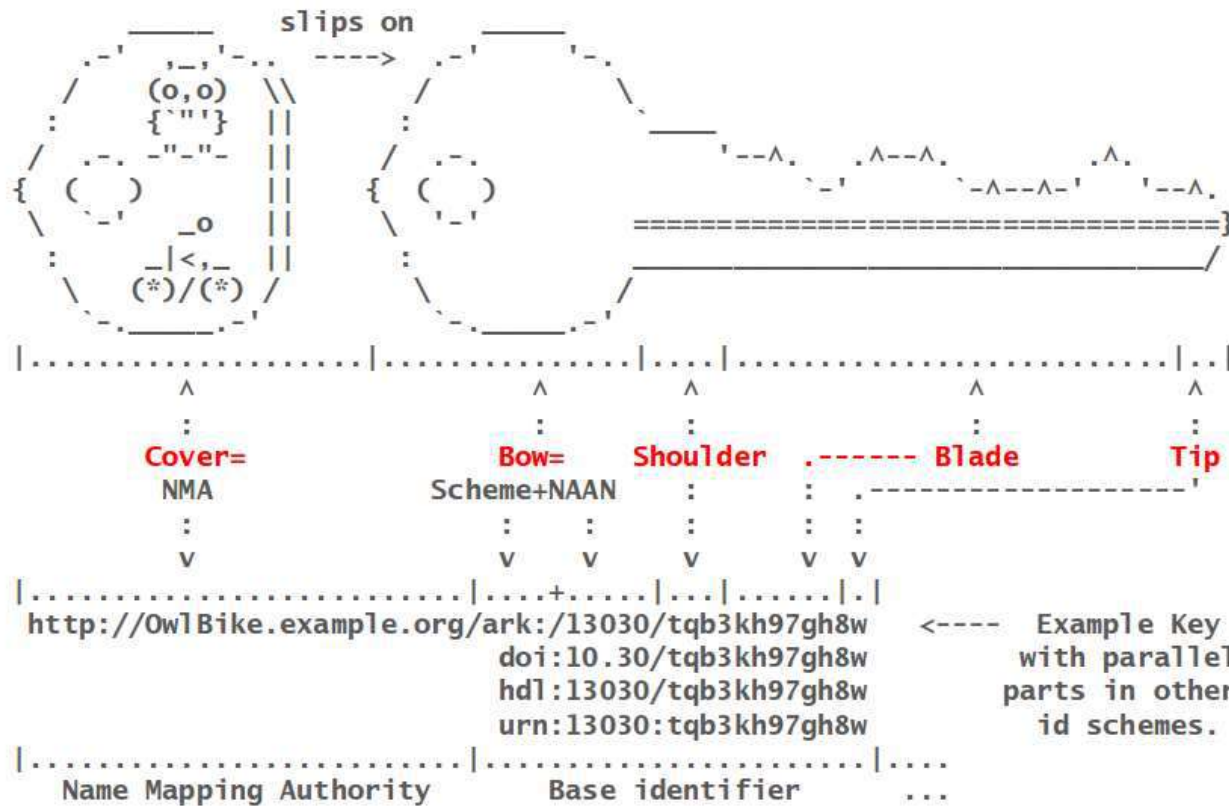
- URLs with long-term support
- Maintained by California Digital Library
- Identify objects of any type (digital, physical, people, vocabulary terms, art...)
- Schema:

`http://example.org/ark:/13030/654xz321/s3/f8.05v.tiff`




ARK - Scheme

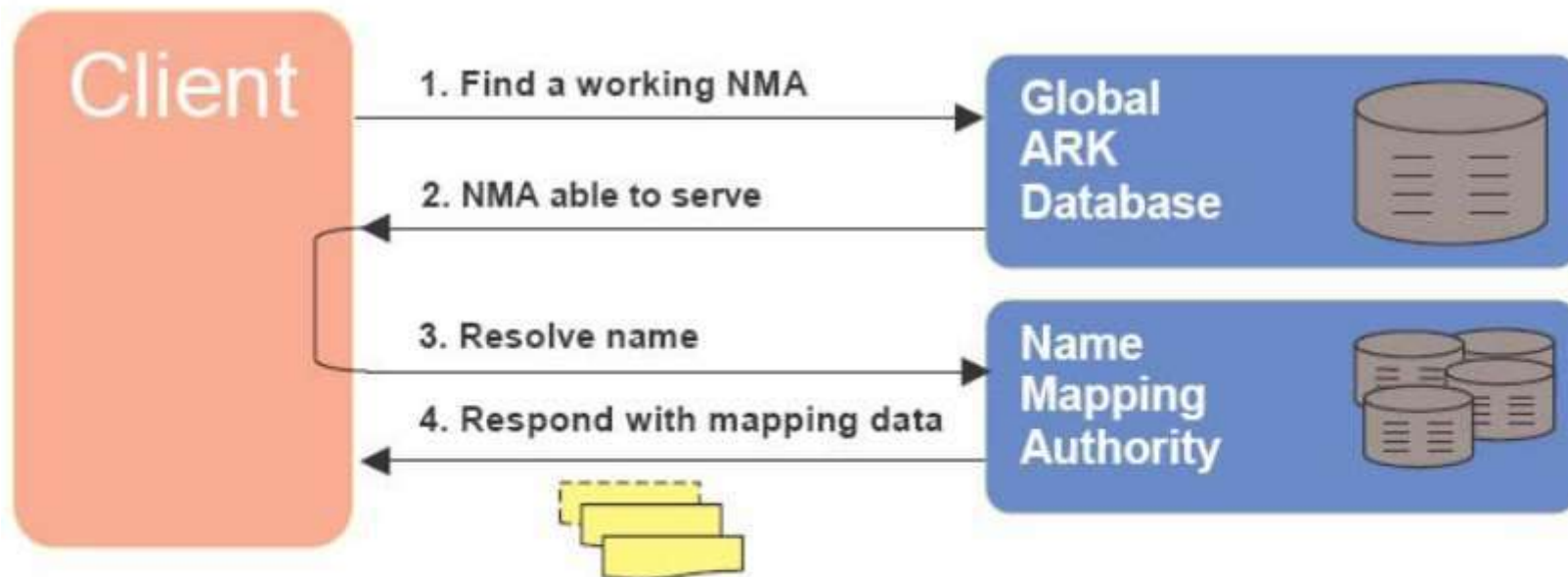
Locksmith jargon: shoulder, blade, tip, bow, cover



ARK

- Currently there are more than 600 NAANs
 - Universities
 - Libraries (e.g. Austrian National Library)
 - Companies (e.g. Google)
 - Organizations (e.g. IEEE) 
 - https://n2t.net/e/pub/naan_table.html
 - http://www.cdlib.org/services/uc3/naan_registry.txt
- Any institution can obtain a NAAN by contacting CDL
- ARK can be self hosted
- ARKs are free

ARK



© MPDL

[10] <http://www.mpi.nl/DAM-LR/meeting5/Persistent%20Identifiers.pdf>

ARK vs. DOI

- ARK
 - Subset facilities
 - Can be deleted
 - Good for early stage of live cycle
 - Free
- DOI
 - **Metadata** cannot be deleted, stored persistently at resolver!
 - Higher reputation
 - Commercial

An Overview of PID Systems

	DOI	ARK	PURL
Actionable	✓	✓	✓
Metadata included	✓	✓	✗
Self hosting	✗	✓	✓
Centralized	✓	✓	✓
Subsets	✓	✓	✗
Opacity	✓	✓	✓
Community Acceptance	✓	✓	✓
Free	✗	✓	✓
Commercial	✓	✗	✗

■ ORCID

- For people (researchers)
- Resolving name ambiguity (common names, name changes)
- Link research activities and output
- ORCID Austria:
<https://www.tuwien.at/kooperationen/orcid/en/home>
- Examples: (increasing order of detailed information provided)
 - <https://orcid.org/0000-0002-9272-6225>
 - <https://orcid.org/0000-0002-4929-7875>

■ Other Person ID systems

- Social Security Number ☺
- Web of Science ResearcherID
- Scopus Author ID



Persistent Identifiers @ TU Wien

■ DOI Service Austria, ORCID-Austria

DOI Service Austria, ORCID Austria

To improve the visibility of Austrian researchers and their academic performance, TU Wien Bibliothek is leading two national initiatives: the DOI Service Austria and ORCID Austria. Not only should this raise awareness of the significance of persistent identifiers (PIDs) in academic communication, it should also create a community of practice.

DOI Service Austria

Since January 2020, with the DOI Service Austria, TU Wien Bibliothek has been providing all Austrian universities, research institutions and other non-profit organisations in the research and education sector domiciled in Austria with an attractive opportunity to register and use Digital Object Identifiers (DOIs) to ensure stable retrieval of academic output via the internet.

For the first time, there is a central point of contact in Austria providing advice on the subject of DOIs, organising international developments and relaying information promptly to Austrian institutions. In order to be able to provide this service as a local authority, TU Wien Bibliothek is a member of the DataCite Association. DataCite is a DOI provider that focuses specifically on the persistent identification of objects stored in repositories and relies on uniform metadata.

Why use DOIs? DOIs are recognised and used internationally. The use of DOIs for research output published on the internet ensures reliable citations and promotes the visibility and stable findability of the document on the internet. The use of DOIs together with other persistent identifiers, such as ORCID iDs for authors and ROR for institutions, also enables improved, reliable and stable attribution of research output to particular persons, research facilities and institutions.

The DOI Service Austria enables Austrian institutions to use the *Fabrica* registration platform and the DataCite interfaces (MDS API, REST API): this enables both manual and automatic registration of DOIs. Customers of the DOI Service Austria receive the prefixes from us for the independent DOI assignment in the respective institutional repositories. TU Wien Bibliothek provides technical support as well as support for the quality assurance of metadata in Austrian information systems and the interoperability between IT applications. Fees for the DOI Service Austria are based on the DataCite cost model. Contact us for more details.

<https://www.tuwien.at/en/library/doi-service-austria-orcid-austria/>

<https://www.tuwien.at/kooperationen/orcid/>
13 Institutions in Austria



Persistent Identifiers for Organizations

- GRID: Global Research Identifier Database
 - <https://www.grid.ac/institutes>
 - grid.5329.d
- ROR: Research Organization Registry
 - <https://ror.org>
 - [04d836q62](https://ror.org/04d836q62)
- Ringgold:
 - <https://www.ringgold.com>
 - RIN 508192

Outline

-
- Why should we want to cite data?
 - What identifier system should I use?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - Summary
-

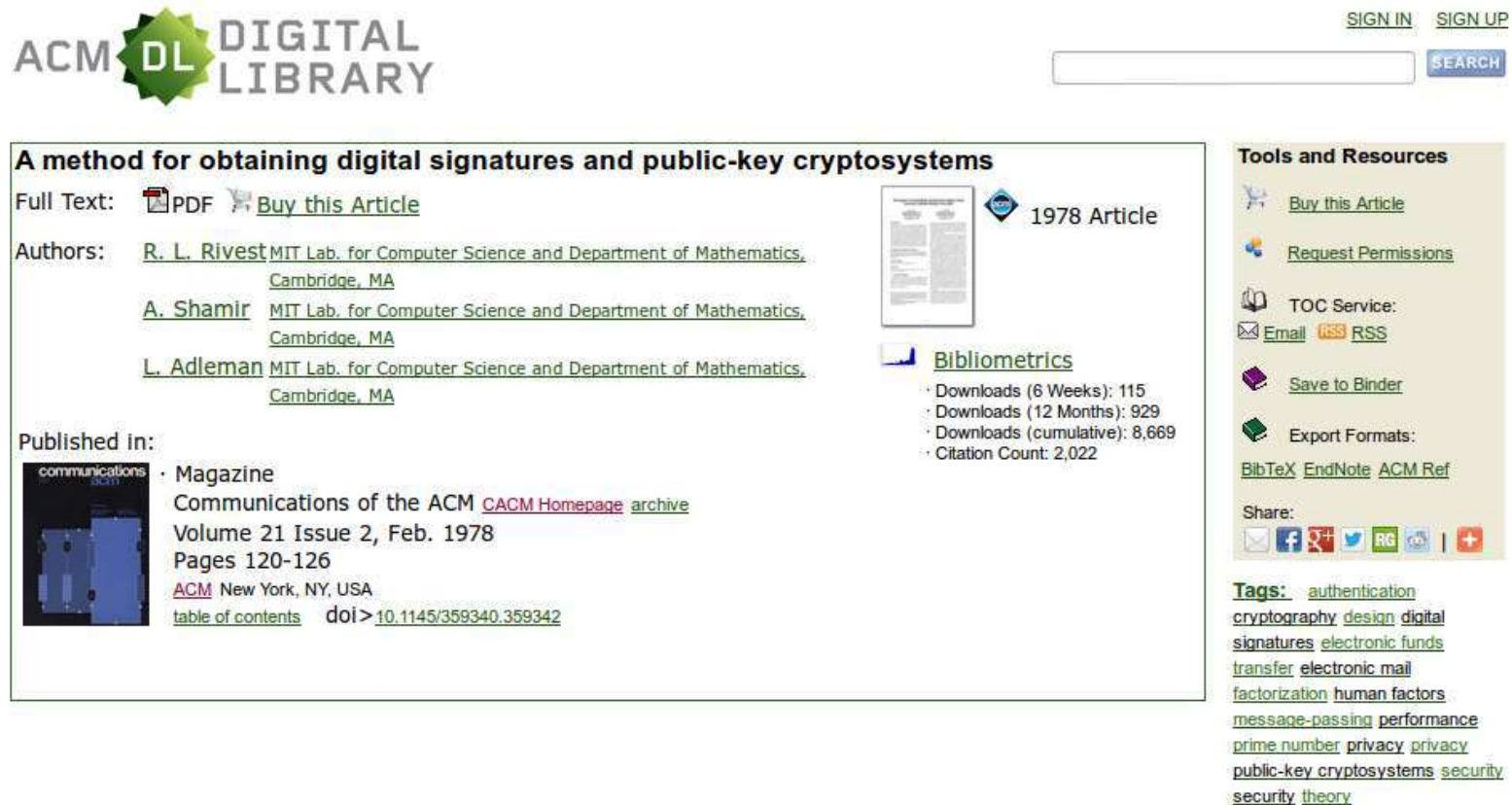
Why to cite data?

- It's what you do! – Lots of benefits
 - Makes live easier because you can build on a solid foundation
 - Speeds up the process because you can re-use existing stuff
 - Helps avoiding / detecting mistakes, improves quality
 - Reuse increases citations, visibility, currency
- But:
 - To achieve this it must be easy, straightforward, “automatic”
 - Citing Papers is easy...
 - ...what about data?
(more about this later... first: “we should just do it”)

Deja-vue

How to cite data?

- Referencing research papers is well established





The screenshot shows the ACM Digital Library interface for the article "A method for obtaining digital signatures and public-key cryptosystems". The page includes a search bar at the top right with "SIGN IN" and "SIGN UP" links. The article title is "A method for obtaining digital signatures and public-key cryptosystems". Below the title, there are links for "Full Text: PDF" and "Buy this Article". The authors listed are R. L. Rivest, A. Shamir, and L. Adleman, all from MIT Lab. for Computer Science and Department of Mathematics, Cambridge, MA. The article was published in "Communications of the ACM" magazine, Volume 21 Issue 2, Feb. 1978, pages 120-126. A "Bibliometrics" section shows download statistics: 115 downloads in 6 weeks, 929 in 12 months, and 8,669 cumulative downloads, with a citation count of 2,022. A "Tools and Resources" sidebar on the right offers options like "Buy this Article", "Request Permissions", "TOC Service", "Email", "RSS", "Save to Binder", and "Export Formats". A "Share:" section includes social media icons for Facebook, Twitter, and others. A "Tags:" section lists various keywords related to the article.

ACM DL DIGITAL LIBRARY


SIGN IN SIGN UP


SEARCH

A method for obtaining digital signatures and public-key cryptosystems

Full Text:  PDF  Buy this Article

Authors: [R. L. Rivest](#) MIT Lab. for Computer Science and Department of Mathematics, Cambridge, MA
[A. Shamir](#) MIT Lab. for Computer Science and Department of Mathematics, Cambridge, MA
[L. Adleman](#) MIT Lab. for Computer Science and Department of Mathematics, Cambridge, MA




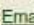

Published in:  Magazine
Communications of the ACM [CACM Homepage](#) [archive](#)
Volume 21 Issue 2, Feb. 1978
Pages 120-126
[ACM](#) New York, NY, USA
[table of contents](#) doi > [10.1145/359340.359342](#)






 1978 Article

Bibliometrics

- Downloads (6 Weeks): 115
- Downloads (12 Months): 929
- Downloads (cumulative): 8,669
- Citation Count: 2,022

Tools and Resources

-  Buy this Article
-  Request Permissions
- TOC Service:
 Email  RSS
-  Save to Binder
- Export Formats:
[BibTeX](#) [EndNote](#) [ACM Ref](#)

Share:     

Tags: [authentication](#) [cryptography](#) [design](#) [digital](#) [signatures](#) [electronic](#) [funds](#) [transfer](#) [electronic mail](#) [factorization](#) [human factors](#) [message-passing](#) [performance](#) [prime number](#) [privacy](#) [privacy](#) [public-key cryptosystems](#) [security](#) [security](#) [theory](#)

- Example: Web-page download

Natural Language Interfaces: What is the Problem? - A data-driven quantitative analysis

Philipp Cimiano¹ and Michael Minock²

¹WIS, TU Delft / ²University of Umea

Abstract. While qualitative analyses of the problems involved in building natural language interfaces (NLIs) have been available, a quantitative grounding in empirical data has been missing. We fill this gap by providing a quantitative analysis on the basis of the Geobase dataset. We hope that this analysis can guide further research in NLIs.

1 Introduction

So far, there has been an impressive amount of research on natural language interfaces (NLIs), i.e. on interfaces allowing users to interact with a certain information system in natural language. While NLIs are not inherently restricted only to the task of answering questions on the basis of a given database or knowledge base, most of the NLIs developed so far have been designed for this purpose. Along these lines, as in most other research on natural language interfaces, we limit ourselves to this restricted view of natural language interfaces essentially as systems providing answers to natural language questions in this paper. Research on NLIs dates back to the 70s and 80s (see [1], [6]) and has yielded increased attention in recent years with a plethora of systems emerging: PRECISE [13], STEP [11], ORAKEL [3], Aqualog [10], GINSENG [2], just to name a few of the very recent systems. What seems missing so far is a description of the problem, in particular a quantitative analysis of the problems inherent in the task of building natural language interfaces. While there have been qualitative analyses of the problems involved in constructing NLIs ([1], [6]), to our knowledge there has been no quantitative analysis grounding the qualitative characteristics of the problem in real data. This is crucial in our view as it can and should guide the development of NLIs in the future, focusing them on the challenging problems. It would also help system developers to focus on a specific phenomenon encountered in NLIs (e.g. resolution of ambiguities) and foster progress in the field by clearly designing and evaluating the solution to a specific phenomenon which would ideally not be specific to one particular approach but reusable across systems. In our view, no real progress can be expected in NLI research only from charts hiding the interesting details and solutions to characteristic problems involved in the task behind top performing precision and recall measures.

The structure of this paper is as follows: in the next Section 2 we describe the dataset we have used to provide a quantitative analysis and describe our methodology. Then, in Section 3 we describe our interesting findings and derive

2 Datasets and Methodology

To provide a quantitative analysis of the problem of constructing an NLI we proceeded as follows: we downloaded a dataset which has been frequently used for the evaluation of natural language interfaces, i.e. the Geobase dataset collected by Mooney and his students¹. The Geobase dataset describes states, cities, mountains, lakes, rivers and roads in the U.S., together with attributes such as area (state, lake), population (state, city), length (river), height (mountain, location) etc.

The datasets consists of a set of 880 test questions (actually 883 questions) and was collected through a web interface hosted at the University of Austin in Texas². We used the 883 test questions for our analysis. After downloading the dataset (in Prolog), we converted the whole dataset into the ontology languages F-Logic [9] and OWL³. The datasets are available from <http://www.cimiano.de> → Projects → Datasets and other Material → ORAKEL.

When converting the dataset into OWL and F-Logic, we used 7 concepts with a total of 17 different relations. We give below the concepts used together with their relations:

Concepts	Relations
state	name, abbreviation, capital, density, population, area, code, hasCity, border, highest_point, lowest_point
city	name, area, inState
river	name, length, flowsThrough
mountain	name, inState, height
road	number, passesThrough
lake	name, area, inState
location	name, inState, height

The design above slightly deviates from the original schema in Mooney's dataset, consisting of 8 relations (state, city, river, border, highlow, mountain, road and lake). We have essentially merged some of the information into one class (the class `state` thus containing the border as well as highest and lowest point information), removed some redundancies (e.g. the name of the state appearing in various relations) and added the `location` class which includes a `height` attribute for the location in question.

The original dataset of Mooney et al. consists of the following 7 relations:

¹ This dataset is available from: <http://www.cs.utexas.edu/users/ml/nldata.html>
² There is also a dataset consisting of 250 questions available from the University of Texas, but this is merely a subset of the larger dataset.
³ <http://www.w3.org/TR/owl-features/>

Example: Web Page Download

To provide a quantitative analysis of the problem of constructing an NLI we proceeded as follows: we downloaded a dataset which has been frequently used for the evaluation of natural language interfaces, i.e. the Geobase dataset collected by Mooney and his students¹. The Geobase dataset describes states, cities, mountains, lakes, rivers and roads in the U.S., together with attributes such as area (state, lake), population (state, city), length (river), height (mountain, location) etc.

The datasets consists of a set of 880 test questions (actually 883 questions) and was collected through a web interface hosted at the University of Austin in Texas². We used the 883 test questions for our analysis. After downloading the dataset (in Prolog), we converted the whole dataset into the ontology languages F-Logic [9] and OWL³. The datasets are available from <http://www.cimiano.de> → Projects → Datasets and other Material → ORAKEL.

¹ This dataset is available from: <http://www.cs.utexas.edu/users/ml/nldata.html>

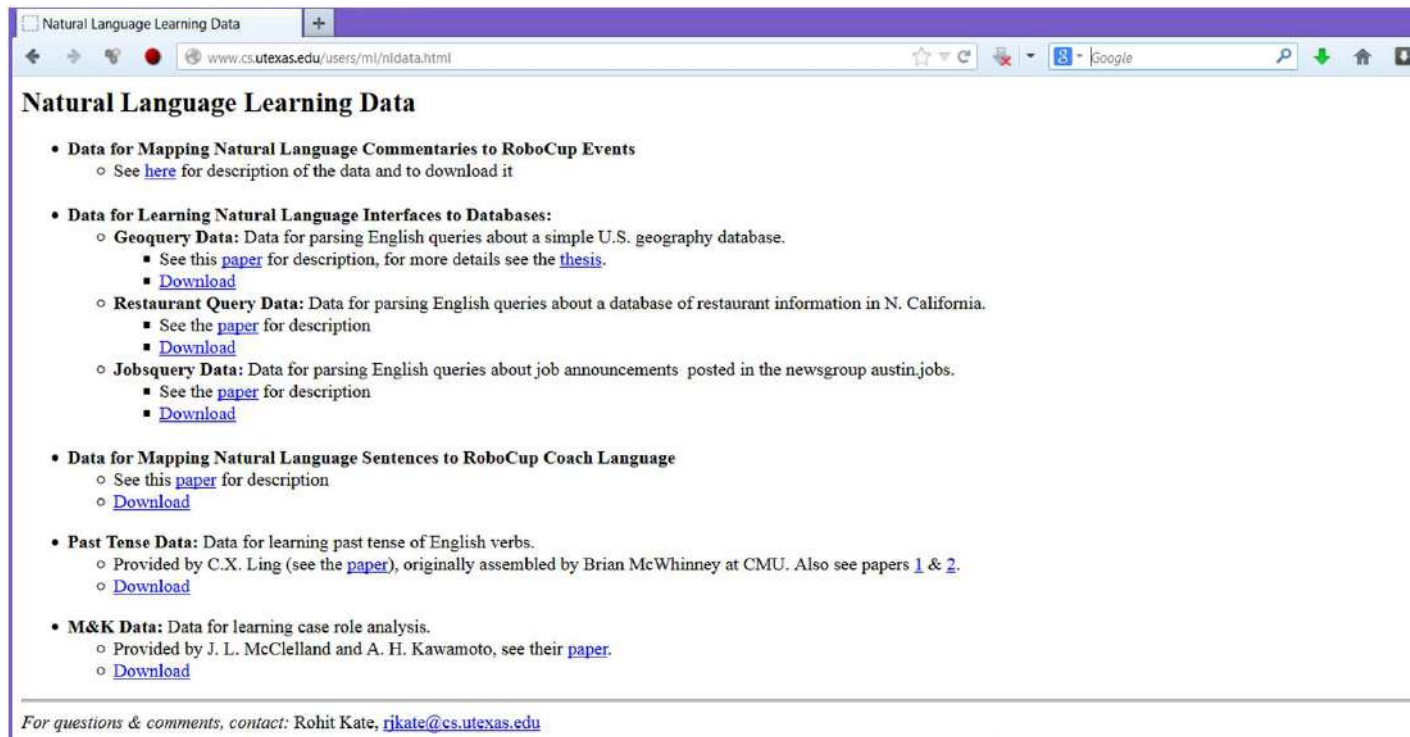
² There is also a dataset consisting of 250 questions available from the University of Texas, but this is merely a subset of the larger dataset.

³ <http://www.w3.org/TR/owl-features/>



Example: Web Page Download

- 1 This dataset is available from: <http://www.cs.utexas.edu/users/ml/nldata.html>
- 2 There is also a dataset consisting of 250 questions available from the University of Texas, but this is merely a subset of the larger dataset.
- 3 <http://www.w3.org/TR/owl-features/>



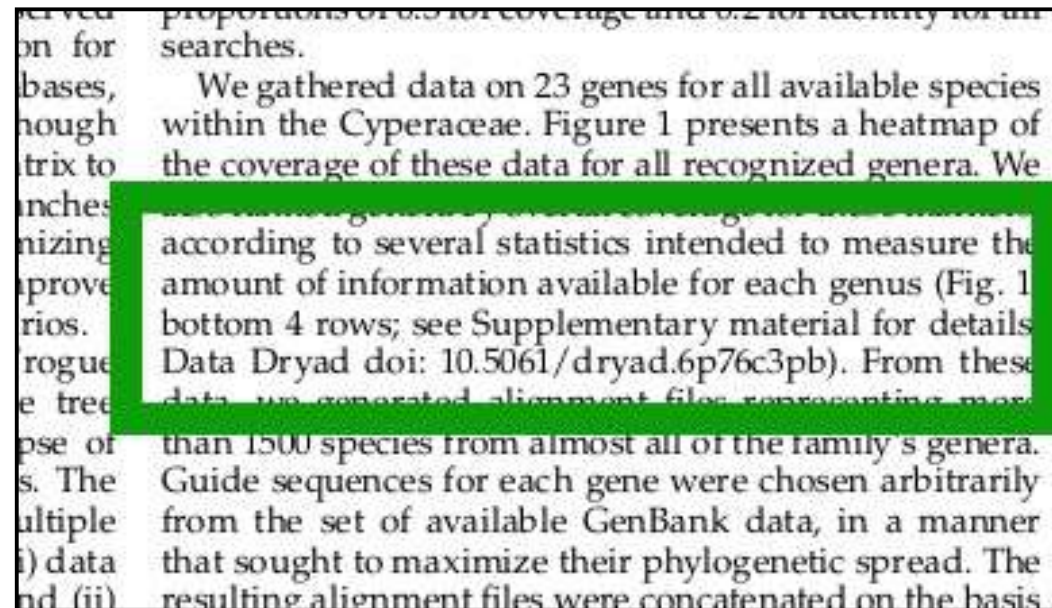
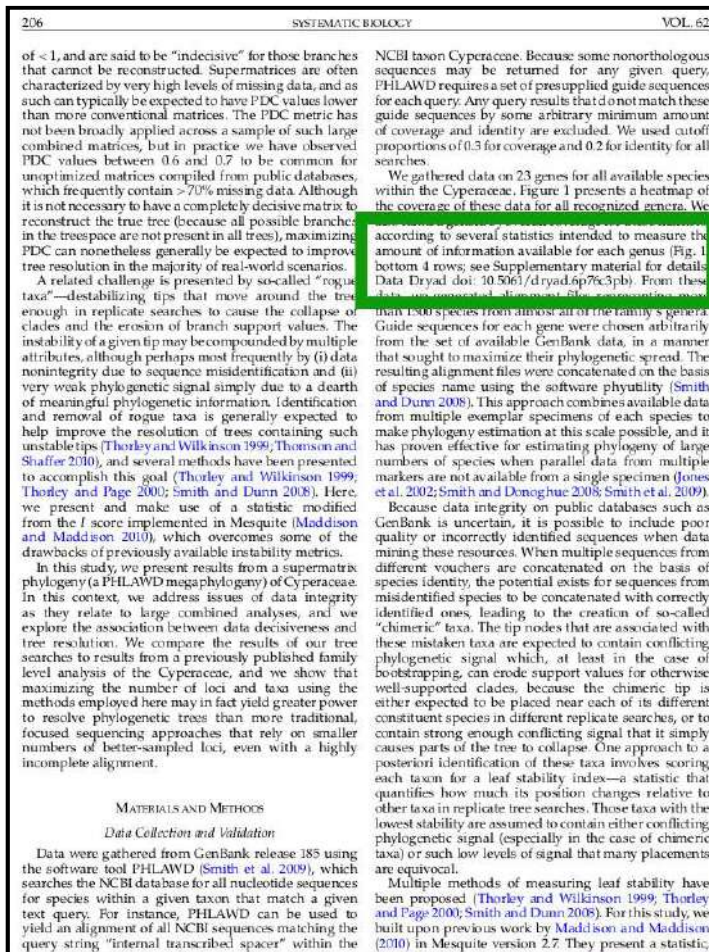
Natural Language Learning Data

- **Data for Mapping Natural Language Commentaries to RoboCup Events**
 - See [here](#) for description of the data and to download it
- **Data for Learning Natural Language Interfaces to Databases:**
 - **Geoquery Data:** Data for parsing English queries about a simple U.S. geography database.
 - See this [paper](#) for description, for more details see the [thesis](#).
 - [Download](#)
 - **Restaurant Query Data:** Data for parsing English queries about a database of restaurant information in N. California.
 - See the [paper](#) for description
 - [Download](#)
 - **Jobsquery Data:** Data for parsing English queries about job announcements posted in the newsgroup austin.jobs.
 - See the [paper](#) for description
 - [Download](#)
- **Data for Mapping Natural Language Sentences to RoboCup Coach Language**
 - See this [paper](#) for description
 - [Download](#)
- **Past Tense Data:** Data for learning past tense of English verbs.
 - Provided by C.X. Ling (see the [paper](#)), originally assembled by Brian McWhinney at CMU. Also see papers [1](#) & [2](#).
 - [Download](#)
- **M&K Data:** Data for learning case role analysis.
 - Provided by J. L. McClelland and A. H. Kawamoto, see their [paper](#).
 - [Download](#)

For questions & comments, contact: Rohit Kate, rjkate@cs.utexas.edu

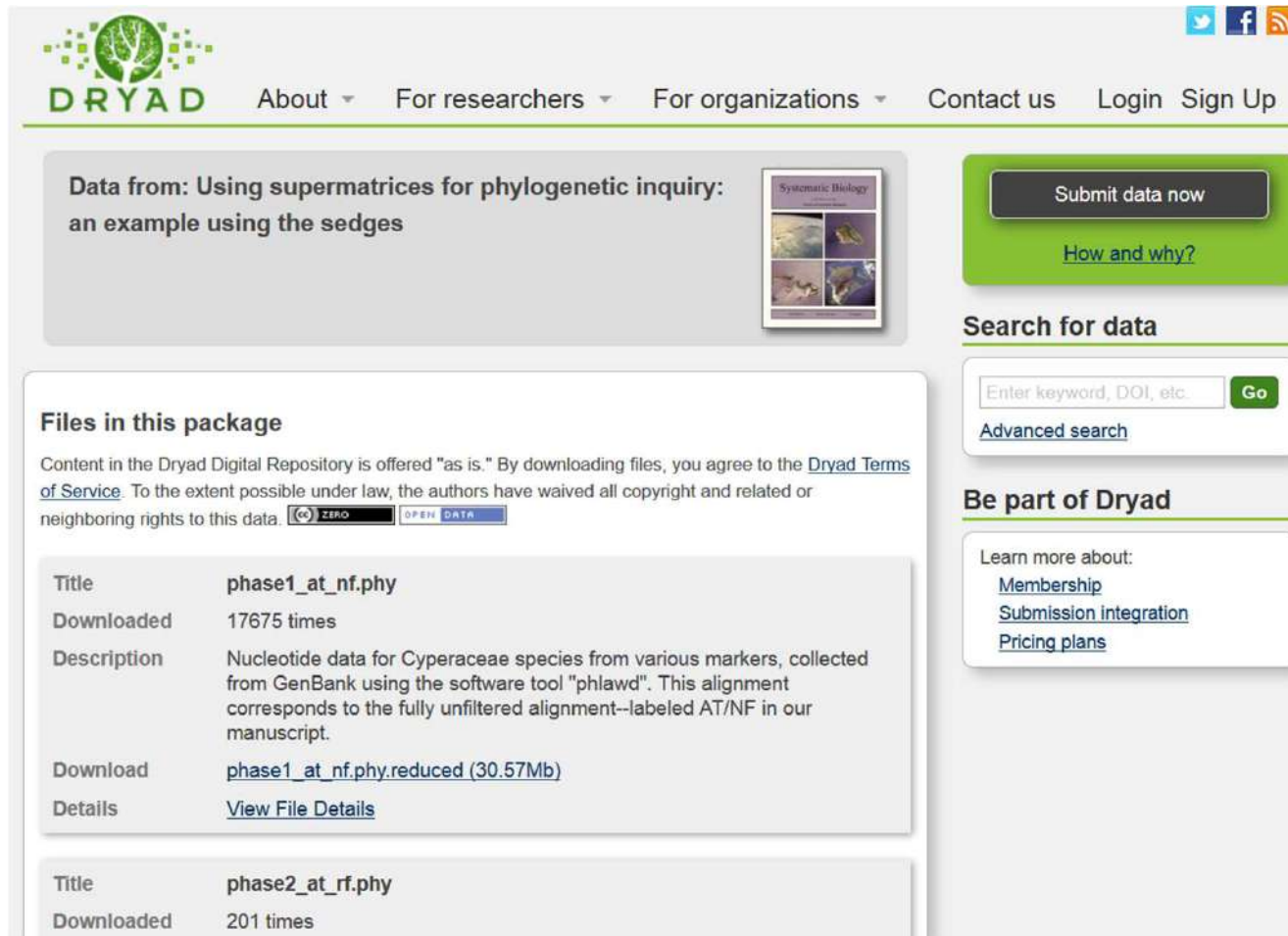
Example: Sharing Platform

- Example: Data sharing platforms



Example: Sharing Platform

- Example: Data sharing platforms



The screenshot shows the Dryad website interface. At the top, there is a navigation bar with the Dryad logo and links for 'About', 'For researchers', 'For organizations', 'Contact us', 'Login', and 'Sign Up'. Below the navigation bar, there is a main content area with a featured article titled 'Data from: Using supermatrices for phylogenetic inquiry: an example using the sedges'. To the right of the article is a 'Submit data now' button and a 'How and why?' link. Below the article is a 'Files in this package' section with a disclaimer and a table of files. To the right of the files section is a search bar and a 'Be part of Dryad' section with links for 'Membership', 'Submission integration', and 'Pricing plans'.

DRYAD About ▾ For researchers ▾ For organizations ▾ Contact us Login Sign Up

Data from: Using supermatrices for phylogenetic inquiry: an example using the sedges

Submit data now
[How and why?](#)


Search for data

Enter keyword, DOI, etc.
[Advanced search](#)

Be part of Dryad

Learn more about:
[Membership](#)
[Submission integration](#)
[Pricing plans](#)

Files in this package

Content in the Dryad Digital Repository is offered "as is." By downloading files, you agree to the [Dryad Terms of Service](#). To the extent possible under law, the authors have waived all copyright and related or neighboring rights to this data. 

Title	phase1_at_nf.phy
Downloaded	17675 times
Description	Nucleotide data for Cyperaceae species from various markers, collected from GenBank using the software tool "phlawd". This alignment corresponds to the fully unfiltered alignment--labeled AT/NF in our manuscript.
Download	phase1_at_nf.phy.reduced (30.57Mb)
Details	View File Details

Title	phase2_at_rf.phy
Downloaded	201 times

<http://datadryad.org/resource/doi:10.5061/dryad.6p76c3pb/1>

Example: Sharing Platform

SYSTEMATIC BIOLOGY VOL. 62

SUPPLEMENTARY MATERIAL

Data files and/or other supplementary information related to this paper have been deposited on Dryad at <http://datadryad.org> under doi: 10.5061/dryad.6p76c3pb.

FUNDING

This work was supported by the National Science

Title	phase2_at_rf.phy
Downloaded	201 times
Description	Nucleotide data for Cyperaceae species from various markers, collected from GenBank using the software tool "phlawd". This alignment corresponds to the unscaffolded, rogues-filtered alignment-labeled AT/RF in our manuscript.
Download	phase2_at_rf.phy.reduced (26.56Mb)
Details	View File Details

Title	phase3_sc_nf.phy
Downloaded	199 times
Description	Nucleotide data for Cyperaceae species from various markers, collected from GenBank using the software tool "phlawd". This alignment corresponds to the scaffolded alignment with rogues unfiltered-labeled SC/NF in our manuscript.
Download	phase3_sc_nf.phy.reduced (7.766Mb)
Details	View File Details

Title	phase4_sc_rf.phy
Downloaded	206 times
Description	Nucleotide data for Cyperaceae species from various markers, collected from GenBank using the software tool "phlawd". This alignment corresponds to the maximally filtered alignment: scaffolded and having had rogues removed-labeled SC/RF in our manuscript.
Download	phase4_sc_rf.phy.reduced (6.976Mb)
Details	View File Details

When using this data, please cite the original publication:

Hinchliff CE, Roalson EH (2012) Using supermatrices for phylogenetic inquiry: an example using the sedges. *Systematic Biology* 62(2): 205-219. <http://dx.doi.org/10.1093/sysbio>

Additionally, please cite the Dryad data package:

Hinchliff CE, Roalson EH (2012) Data from: Using supermatrices for phylogenetic inquiry: an example using the sedges. Dryad Digital Repository. <http://dx.doi.org/10.5061/dryad.6p76c3pb>

[Cite](#) | [Share](#)

<http://datadryad.org/resource/doi:10.5061/dryad.6p76c3pb/1>

- Example: Subsets of data



Available online at www.sciencedirect.com



Pattern Recognition Letters 28 (2007) 1133–1141

Pattern Recognition Letters

www.elsevier.com/locate/patrec

Introducing a very large dataset of handwritten Farsi digits and a study on their varieties

Hossein Khosravi^{a,b,*}, Ehsanollah Kabir^a

^a *Department of Electrical Engineering, Tarbiat Modarres University, Tehran, Iran*
^b *Research and Development Unit, HODA System Co., Tehran, Iran*

Received 5 September 2005; received in revised form 24 September 2006
Available online 21 February 2007

Communicated by A. M. Alimi

Abstract

A very large dataset of handwritten Farsi digits is introduced. Binary images of 102,352 digits were extracted from about 12,000 registration forms of two types, filled by B.Sc. and senior high school students. These forms were scanned at 200 dpi with a high speed scanner. A method for finding variety of handwritten digits in a typical dataset is proposed. Based on this method, training and test subset are provided to facilitate sharing of results among researchers as well as performance comparison.

© 2007 Elsevier B.V. All rights reserved.

- Example: Subsets of data

1134 *H. Khosravi, E. Kabir / Pattern Recognition Letters 28 (2007) 1133–1141*

Table 1
Some popular digit datasets

Dataset	dpi	Training samples	Test samples	Total samples
CENPARMI	166	4000	2000	6000
CEDAR	300	18,468	2711	21,179
MNIST	Normalized into 28 * 28	60,000	10,000	70,000
USPS	300	7291	2007	9298

The CEDAR³ digit dataset is available from CEDAR, SUNY⁴ at Buffalo. The images were scanned at 300 dpi. The training and test sets contain 18468 and 2711 digits, respectively. The number of samples in both training and test sets differ for each class. Since some images in the test set are poorly segmented, a subset of 2213 well-segmented images are also provided for testing (Liu, 2003).

The MNIST, modified NIST⁵ dataset (LeCun et al., 1995) was extracted from the NIST datasets SD3 and SD7. The training and test sets are composed from both SD3 and SD7. Samples are normalized into 28 * 28 gray-scale images with aspect ratio reserved, and the normalized images are located in a 28 * 28 frame. The dataset is available from LeCun. Number of training and test samples are 60,000 and 10,000 respectively.

At last the USPS digit dataset has 7291 training and 2007 test samples (Hull, 1994). Table 1 lists these datasets briefly.




Fig. 1. Sample handwritten Farsi digits.

There were several fields in both types of forms. We used two digit fields from type 1, including *Postal Code* and *National Code*, each of 10 digits length and three digit fields from type 2 including *Record Number*, *Identity Certificate Number* and *Phone Number* that at most have 26 digits, while in average about 20 digits. Both forms are in color. In both types, handwritten texts are in blue or occasionally in black.

3.2. Digit extraction and recognition

To extract the digits, we must find the regions of interest. There were at least two reference marks (squares) in each form (circled in Fig. 2). We first search for these marks using a simple and fast algorithm shown in Fig. 3. If they are not found, the form is rejected. This situation occurs rarely, e.g. when the paper is scanned upside down or the reference square is too noisy. Then, if the reference squares are not in their expected positions, the form is rotated and shifted so that these squares are placed in the

- Example: Subsets of data

5. Choosing the training and test sets

To facilitate sharing of results on this dataset between researchers, we provide two distinct datasets for training and test.

From Table 3 it can be seen that the most usual styles are fallen into samples S1, and other varieties are fallen into S2, S3 and S4. So we tried to select most of training samples from S1. To be more accurate we selected from each category a number of samples equal to their proportion in total samples, i.e. 73.47% of training samples were selected from S1, 9.83% from S2 and so on. Then we set aside training samples and select test samples from the remaining samples, randomly. In this way the training set is a true representation of the whole population, while the test set is selected without any predefined information.

We selected 60,000 samples for training set and 20,000 for test. The remaining samples are also available in another subset (see Appendix A).

Appendix A. Dataset specification and availability

The dataset is available in four separate files, **Total.cdb**, **Training.cdb**, **Test.cdb**, **Remaining.cdb**. The file format is described here with a pseudo code:

```

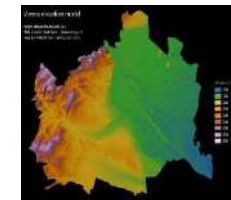
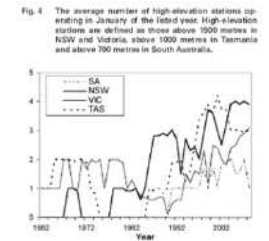
Skip Header (1024 bytes)
while not End of File
{
    read Start Byte: (1 byte) 0xFF that
    specifies the start of new image
    read Label: (1 byte) character label
    read Width: (1 byte) character width
    read Height: (1 byte) character height
    read Byte Count: (2 bytes) number of bytes for this character.
    //Runlength coding on each row
    for y=0 to Height
        while(x < Width)
        {
            read NumOfWhitePixels,
            read NumOfBlackPixels;
        }
    }
}
    
```

Source codes for reading the dataset files are available in Matlab, C++ and Pascal. To get the dataset please contact kabir@modares.ac.ir, or see the homepage <http://www.modares.ac.ir/eng/kabir>.



Motivation

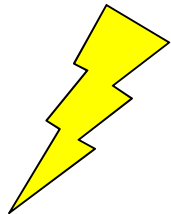
- Research data is fundamental for science/industry/...
 - Data serves as input for workflows and experiments
 - Data is the source for graphs and visualisations in publications
 - Decisions are based on data
- Data is needed for Reproducibility
 - Repeat experiments
 - Verify / compare results
- Need to provide specific data set
 - Service for data repositories



<https://commons.wikimedia.org/w/index.php?curid=30978545>

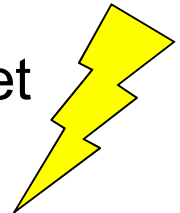
1. Put data in data repository,
2. Assign PID (DOI, Ark, URI, ...)
3. Make is accessible
→ done!?

Identification of Dynamic Data

- Usually, datasets have to be static
 - Fixed set of data, no changes:
no corrections to errors, no new data being added
 - But: (research) data is **dynamic**
 - Adding new data, correcting errors, enhancing data quality, ...
 - Changes sometimes highly dynamic, at irregular intervals
 - Current approaches
 - Identifying entire data stream, without any versioning
 - Using “accessed at” date
 - “Artificial” versioning by identifying batches of data (e.g. annual), aggregating changes into releases (time-delayed!)
- 
- Would like to identify precisely the **data as it existed at a specific point in time**

Granularity of Subsets

- What about the **granularity** of data to be identified?
 - Enormous amounts of CSV data
 - Researchers use specific subsets of data
 - Need to identify precisely the subset used
 - Current approaches
 - Storing a copy of subset as used in study -> scalability
 - Citing entire dataset, providing textual description of subset -> imprecise (ambiguity)
 - Storing list of record identifiers in subset -> scalability, not for arbitrary subsets (e.g. when not entire record selected)
- Would like to be able to identify precisely the **subset of (dynamic) data used** in a process



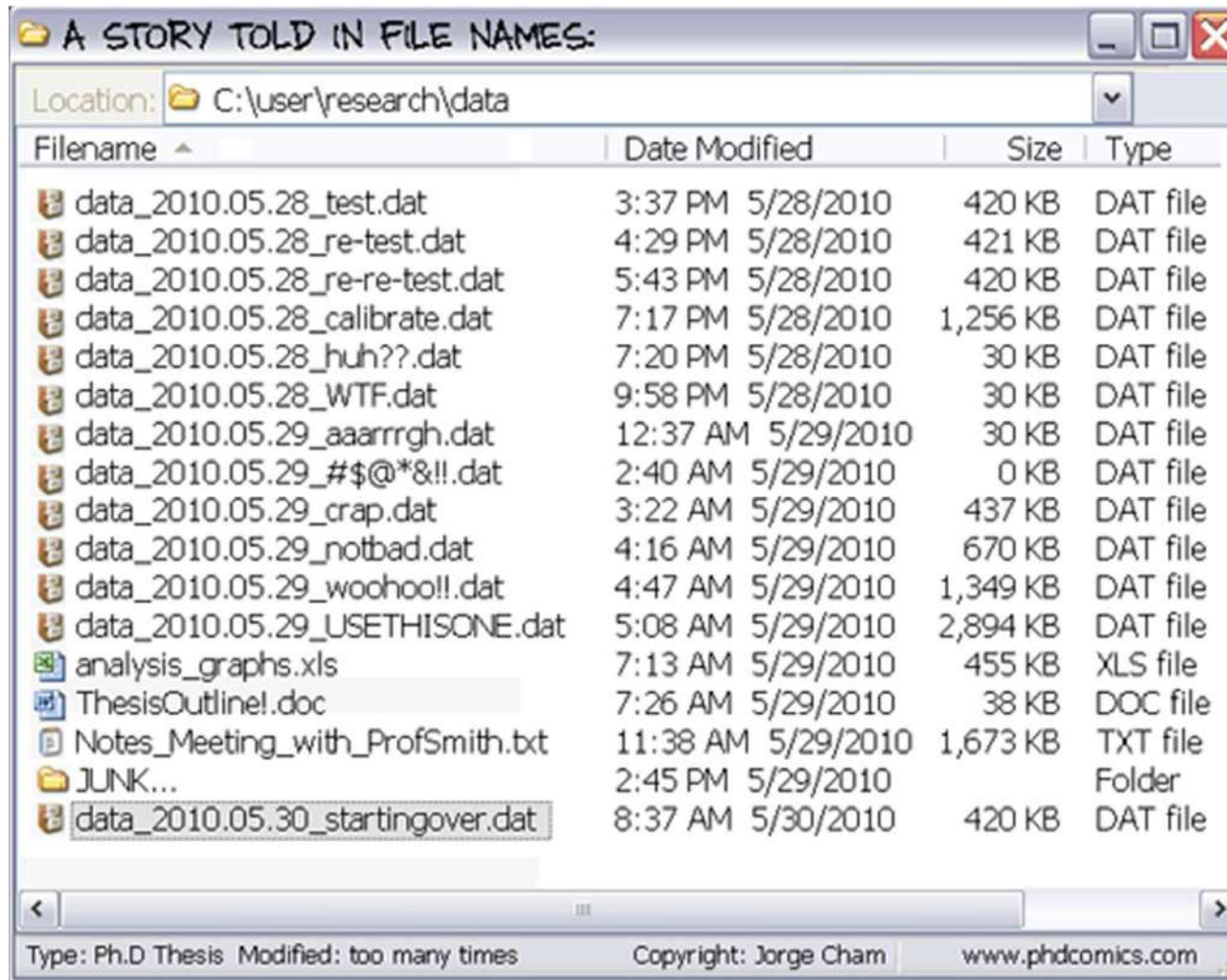
Data Citation – Requirements

- Dynamic data
 - corrections, additions, ...
- Arbitrary subsets of data (granularity)
 - rows/columns, time sequences, ...
 - from single number to the entire set
- Stable across technology changes
 - e.g. migration to new database
- Machine-actionable
 - not just machine-readable,
definitely not just human-readable and interpretable
- Scalable to very large / highly dynamic datasets
 - But: should also work for small and/or static datasets!



What we do NOT want...

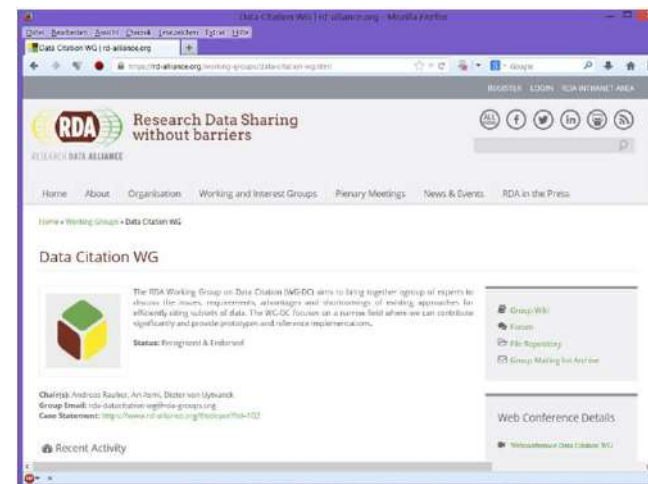
- Common approaches to data management...
(from PhD Comics: A Story Told in File Names, 28.5.2010)
Source: <http://www.phdcomics.com/comics.php?f=1323>



Outline

-
- Why should we want to cite data?
 - What identifier system should I use?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - Summary
-

- Research Data Alliance
- WG on **Data Citation: Making Dynamic Data Citeable**
- March 2014 – September 2015
 - Concentrating on the problems of **large, dynamic (changing) datasets**
- Final version presented Sep 2015 at P7 in Paris, France
- Endorsed September 2016 at P8 in Denver, CO
- Since then: supporting adopters



<https://www.rd-alliance.org/groups/data-citation-wg.html>



RDA WGDC - Solution

- **We have**
 - Data &  some means of access („query“)



Dynamic Data Citation

We have: Data + Means-of-access

**Dynamic Data Citation:
Cite (dynamic) data dynamically via query!**

Dynamic Data Citation

We have: Data + Means-of-access

**Dynamic Data Citation:
Cite (dynamic) data dynamically via query!**

Steps:

1. Data → versioned (history, with time-stamps) 

We have: Data + Means-of-access

**Dynamic Data Citation:
Cite (dynamic) data dynamically via query!**

Steps:

1. Data → versioned (history, with time-stamps)

Researcher creates working-set via some interface:

We have: Data + Means-of-access

**Dynamic Data Citation:
Cite (dynamic) data dynamically via query!**

Steps:

1. Data → versioned (history, with time-stamps)

Researcher creates working-set via some interface:

2. Access → **store & assign PID to “QUERY”**, enhanced with
- **Time-stamping** for re-execution against versioned DB
 - **Re-writing** for normalization, unique-sort, mapping to history
 - **Hashing** result-set: verifying identity/correctness

leading to landing page

S. Pröll, A. Rauber. **Scalable Data Citation in Dynamic Large Databases: Model and Reference Implementation**. In IEEE Intl. Conf. on Big Data 2013 (IEEE BigData2013), 2013

http://www.ifs.tuwien.ac.at/~andi/publications/pdf/pro_ieeebigdata13.pdf

Data Citation – Deployment

- Researcher uses workbench to identify subset of data
- Upon executing selection („download“) user gets
 - Data (package, access API, ...)
 - PID (e.g. DOI) (Query is time-stamped and stored)
 - Hash value computed over the data for local storage
 - Recommended citation text (e.g. BibTeX)
- PID resolves to landing page
 - Provides detailed metadata, link to parent data set, subset,...
 - Option to retrieve original data OR current version OR changes
- Upon activating PID associated with a data citation
 - Query is re-executed against time-stamped and versioned DB
 - Results as above are returned
- Query store aggregates data usage



Data Citation – Deployment

- Note: query string provides excellent subset of data
- provenance information on the data set! user gets
 - Data (package, access API, ...)
 - PID (e.g. DOI) (Query is time-stamped and stored)
 - Hash value computed over the data for local storage
 - Recommended citation text (e.g. BibTeX)
- PID resolves to landing page
 - Provides detailed metadata, link to parent data set, subset,...
 - Option to retrieve original data OR current version OR changes
- Upon activating PID associated with a data citation
 - Query is re-executed against time-stamped and versioned DB
 - Results as above are returned
- Query store aggregates data usage

Data Citation – Deployment

- Note: query string provides excellent subset of data
- provenance information on the data set! user gets
 - Data (package)
 - PID (e.g. [10.616/1123456789](#))
 - Hash value
 - Recommended citation text (e.g. PID TEX)
- PID resolves to landing page
 - Provides detailed metadata, link to parent data set, subset,...
 - Option to retrieve original data OR current version OR changes
- Upon activating PID associated with a data citation
 - Query is re-executed against time-stamped and versioned DB
 - Results as above are returned
- Query store aggregates data usage

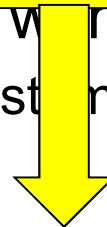
This is an important advantage over traditional approaches relying on, e.g. storing a list of identifiers/DB dump!!!

Data Citation – Deployment

- Note: query string provides excellent subset of data
- provenance information on the data set! user gets
 - Data (package)
 - PID (e.g. DOI)
 - Hash value
 - Recommended citation text (e.g. PID text)
- PID resolves
 - Provides details
 - Option to retrieve
- Upon activating PID associated with a data citation
 - Query is re-executed against time-stamped and versioned DB
 - Results as above are returned
- Query store aggregates data usage

This is an important advantage over traditional approaches relying on, e.g. storing a list of identifiers/DB dump!!!

Identify which parts of the data are used. If data changes, identify which queries (studies) are affected



Data Citation – Output

- 14 Recommendations grouped into 4 phases:
 - 2-page flyer <https://rd-alliance.org/recommendations-working-group-data-citation-revision-oct-20-2015.html>
- Detailed report: Bulletin of IEEE TCDL 2016 http://www.ieee-tcdl.org/Bulletin/v12n1/papers/IEEE-TCDL-DC-2016_paper_1.pdf
- Adopter's reports, webinars <https://www.rd-alliance.org/group/data-citation-wg/webconference-data-citation-wg.html>
- Review / Lessons Learned
Andreas Rauber et al., Precisely and Persistently Identifying and Citing Arbitrary Subsets of Dynamic Data
- Harvard Data Science Review, 3(4), 2021.
DOI [10.1162/99608f92.be565013](https://doi.org/10.1162/99608f92.be565013).



Data Citation – Recommendations

Preparing Data & Query Store

- R1 – Data Versioning
- R2 – Timestamping
- R3 – Query Store

When Resolving a PID

- R11 – Landing Page
- R12 – Machine Actionability

When Data should be persisted

- R4 – Query Uniqueness
- R5 – Stable Sorting
- R6 – Result Set Verification
- R7 – Query Timestamping
- R8 – Query PID
- R9 – Store Query
- R10 – Citation Text

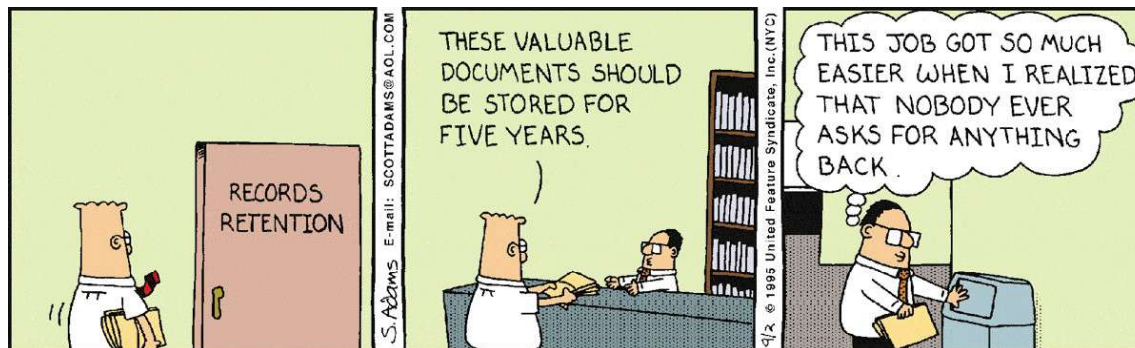
Upon Modifications to the Data Infrastructure

- R13 – Technology Migration
- R14 – Migration Verification



R1: Data Versioning

- **Apply versioning to ensure earlier states of the data can be retrieved**
- Versioning allows tracing the changes (static data: no changes – principle still applies)
- No in-place updates or deletes
 - Mark record as deleted, re-insert new record instead of update
 - Keep old versions – only way to be able to “go back”
- Do we really need to keep everything?
 - (*“changes that were never read never existed”*)



Src: <http://dilbert.com/strip/1995-09-02>

R2: Data Timestamping


- **Ensure that operations on data are timestamped, i.e. any additions, deletions are marked with a timestamp**
- Timestamping is closely related to versioning
- Granularity depends on
 - Change frequency / tracking requirements
 - Per individual operation
 - Batch-operations
 - Grouped in-between read accesses (*“changes that were never read do not matter”*)
 - System (data storage, databases)
 - e.g. FAT 2 seconds, NTFS 100 ns, EXT4 1 ns



https://www-03.ibm.com/ibm/history/exhibits/cc/cc_T30.html

R1 & R2: Versioning / Timestamping

Note:

- R1 & R2 are already pretty much standard in many (RDBMS-) research databases
- Different ways to implement, depending on
 - data type / data structure: RDBMS, CSV, XML, LOD, ...
 - data volume
 - amount and type of changes
 - number of APIs, flexibility to change them
- Distributed settings:
 - synchronized clocks, or: 
 - each node keeps individual, local time
time-stamps for distributed queries based on local times
these local times are stored at the query store aggregating the results



Timestamping vs. Semantic Versioning

Why timestamps, why not semantic versioning

- Some prefer to use semantic versioning (minor/major updates that do not / do change behaviour/interface)
 - Advantage: version number indicates relationship btw. versions
 - Disadvantage:
 - Something that was expected to be a not-changing update may turn out to induce changes / side-effects later-on
 - With data, “minor” updates are hard to think of: changing a typo may result in a record being found / not found by a query, encoding changes may break subsequent processing pipelines
 - Different semantics / types of use across different communities
- Recommendation
 - No semantics in identifier (mantra!)
 - Keep identification (version timestamp) and semantics separate
 - Semantic version number in addition to timestamp



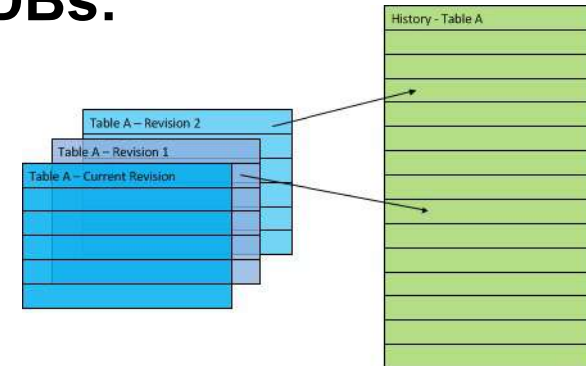
Semantic Versioning

- Semantic versions are “only” **assertions on states of the data at certain points in time**, eg
 - Data may be transient / still undergoing changes, whereas after a certain points in time it has reached a state where no further changes are expected
 - Certain states of data may not be intended for permanent retention, whereas others may have guarantees of availability over time
- Assertions specified as tags associated to queries, e.g.
 - Query *“Select * FROM <table> WHERE timestamp_added < ts1 and ts_deleted >ts1”* may carry the assertions *“status: not expected to change”* and *“availability: 7 years”* (preferably from controlled vocabularies)
- *Subset queries are “nested queries” on such “stable versions”*

R1 & R2: Versioning / Timestamping

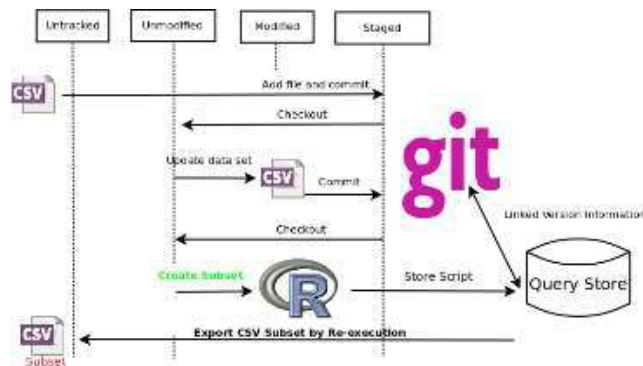
Implementation options for e.g. relational DBs:

- History Table
 - Utilizes full history table
 - Also inserts reflected in history table
 - Doubles storage space, no API adaptations
- Integrated
 - Extend original tables by temporal metadata
 - Expand primary key by timestamp/version column
 - Minimal storage footprint, changes to all APIs
- Hybrid
 - Utilize history table for deleted record versions with metadata
 - Original table reflects latest version only
 - Minimal storage footprint, some API change, expensive query re-writes
- Solution to be adopted depends on trade-off
 - Storage Demand
 - Query Complexity
 - Software/API adaption



Git Implementation 1

- Upload CSV files to Git repository (versioning)
- Subsets created via scripting language (e.g. R)
 - Select rows/columns, sort, returns CSV + metadata file
 - Metadata file with script parameters stored in Git
 - (Scripts stored in Git as well)
- PID assigned to metadata file
 - Use Git to retrieve proper data set version and re-execute script on retrieved file



```

diff --git a/suppList.csv b/suppList.csv
index fe2799f..986025 100644
--- a/suppList.csv
+++ b/suppList.csv
@@ -2,23 +1,21 @@
#suppnameNumber,work,doi,urlNumber,versionInteger,mail
0,C201204,0,204248298228862742,228310,johfey.wilkinson@gmail.com
1,com_sov,0,204629619717021,737174,netc.welsh@epfl.ch
2,Dai_rasm,0,4125818928205009,738220,sapar.conroy@gmail.com
3,ori_gin,0,204809812528468,888358,arsh_ayyoubi@gmail.com
4,Orma_ni,0,20274562734612136,20850,marc_baeupf@gmail.com
5,laro_occ,0,2244266822428982,222170,joell.predator@gmail.com
6,Reynolds,0,15567703313425617,833678,nreynold@earth.umd.edu
7,Glac_sail,0,543290284312615,121550,marco.lajp@epfl.ch
8,2015_708,0,12134819492304000,293452,293452@cs.stonybrook.edu
9,Tatem_fu,0,2415664117016023,202752,Horino.yuki@bunkyo.ac.jp
10,Ipsam_qa,0,2025800473120092,204259,logan.lens@yale.edu
11,Santoro,0,2026801247421,184256,corbyn_howe@gmail.com
12,F4_coll,0,2025800473120092,204259,logan.lens@yale.edu
13,et_cc,0,20264049468228794,248625,johndoherty2550@gmail.com
14,Carey_d,0,221103882266707,220110,joel.welsh@epfl.ch
15,ml_olag,0,2245868222128229,221849,joel.welsh@epfl.ch
16,296684448892,221870,joell.predator@gmail.com
17,0201020229134,221871,joel.welsh@epfl.ch
18,42122324234234,423423,joel.welsh@epfl.ch
19,047042424234234,423423,joel.welsh@epfl.ch
20,222072424234234,222072,joel.welsh@epfl.ch
21,353092424234234,353092,joel.welsh@epfl.ch
22,0201020229134,221870,joel.welsh@epfl.ch
23,346682424234234,346682,joel.welsh@epfl.ch
24,525092424234234,525092,joel.welsh@epfl.ch
25,106012424234234,106012,joel.welsh@epfl.ch
26,420092424234234,420092,joel.welsh@epfl.ch
27,675892424234234,675892,joel.welsh@epfl.ch
28,2042012424234234,204201,joel.welsh@epfl.ch
29,33242092424234234,332420,joel.welsh@epfl.ch
30,1656171716023,202752,Horino.yuki@bunkyo.ac.jp
31,2025732281568,202210,corbyn_howe@gmail.com
32,2020831342617,833678,nreynold@earth.umd.edu
  
```

```

# PID=1234/abcd5678
# Repository: Path=/media/Data/Git-Repository
# Execution_Time=2015-09-30:11:07:09
# Subset_Tool=scripting_front-end version 3.0.2 (2015-08-14)
# Subset_Tool_Path=/usr/bin/Rscript
# Input_Script_Path=/supercomputing/top5-script.v
# Input_Script_Hash=bf3d...d7881:supercomputing/top5-script.r
# Dataset_Path=/supercomputing/supercomputer.csv
# Dataset_Commit_Hash=c9cd...d78c:supercomputer.csv
# Output_Path=/tmp/supercomputer-top5.csv

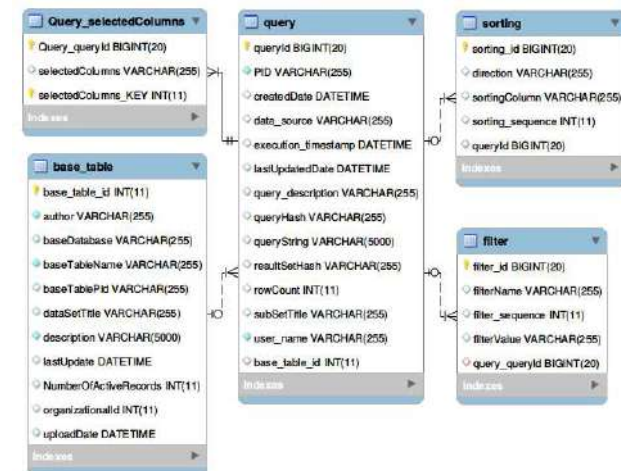
# Original execution:
# /usr/bin/Rscript /supercomputing/top5-script.v \
# /media/Data/Git-Repository/supercomputing/supercomputer.csv \
# /tmp/supercomputer-top5.csv

# Recommended re-execution
# Retrieve script
git --git-dir=/media/Data/Git-Repository/.git/ \
show bf3d...d7881:supercomputing/top5-script.r
> /tmp/reproduced-datasets/top5-script.r
# Retrieve data set
git --git-dir=/media/Data/Git-Repository/.git/ \
show d78d...b9792:supercomputing/supercomputer.csv \
> /tmp/reproduced-datasets/supercomputer.csv

# Re-execute
/usr/bin/Rscript /supercomputing/top5-script.r \
/tmp/reproduced-datasets/supercomputer.csv \
/tmp/reproduced-datasets/supercomputer-top5.csv
  
```

R3: Query Store

- Provide means for storing queries and the associated metadata in order to re-execute them.
- Approach is based upon queries.
 - Therefore we need to preserve the queries
 - Original and re-written (**R4, R5**), potentially migrated (**R13**)
 - Query parameters and system settings
 - Execution metadata
 - Hash keys (multiple, if re-written) (**R4, R6**)
 - **Persistent identifier(s)** (**R8**)
 - Citation text (**R10**) ...
- Comparatively small, even for high query volumes



R4: Query Uniqueness

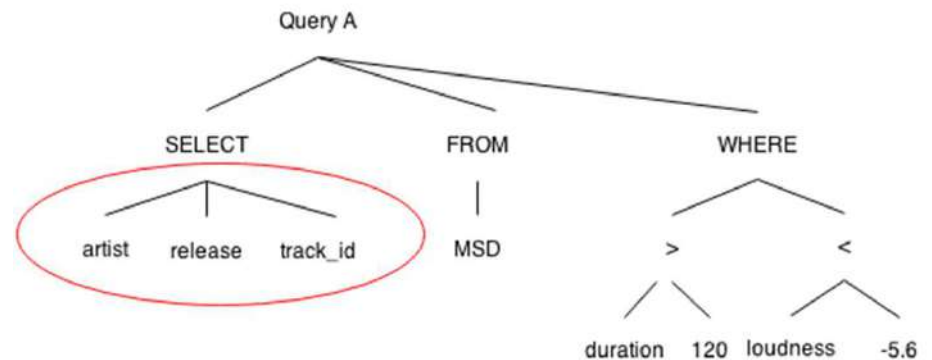
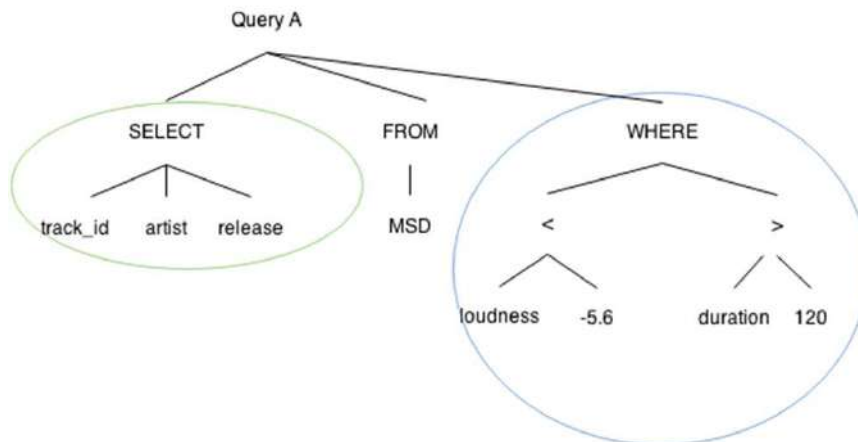
- **Re-write the query to a normalized form so that identical queries can be detected.
Compute checksum of the normalized query to efficiently detect identical queries**
- Detecting identical queries can be challenging
 - Query semantics can be expressed in different ways
 - Different queries can deliver identical results
 - Interfaces can be used for maintaining a stable query structure
- Best effort, no perfect solution
- Usually not a problem if queries generated via standardized interfaces, e.g. workbench – optional!
- Worst case: two PIDs for semantically equivalent queries

R4: Query Uniqueness

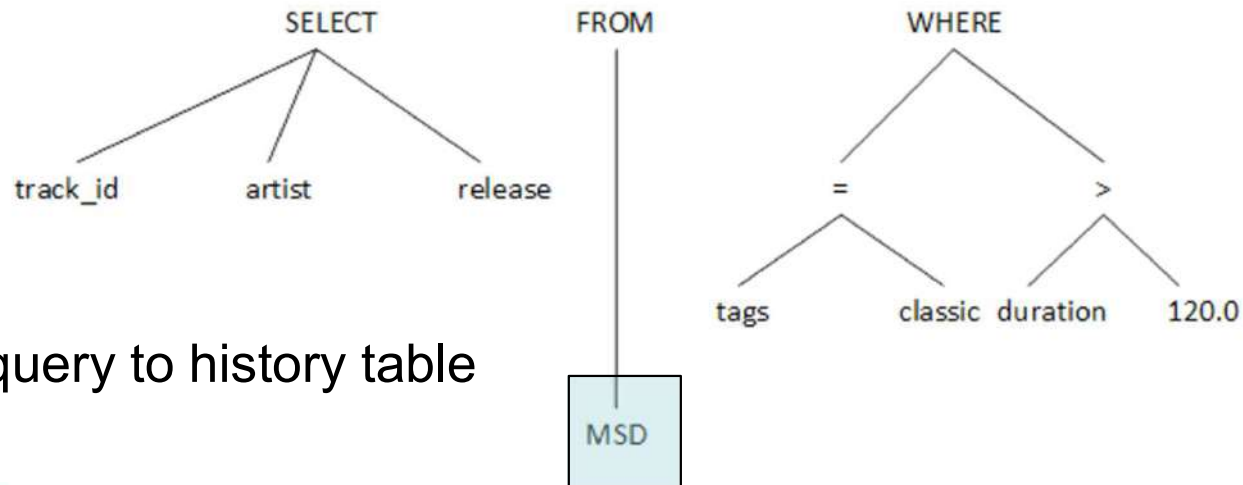
- Query re-writing needed to
 - **Standardization/Normalization** of query to help with identifying semantically identical queries
 - upper/lower case spelling, sorting of filter parameters, ...
 - Re-write to **adapt to versioning approach** chosen (versioning in operational tables, separate history table, ...), e.g. **identify last change to result set touched** upon (i.e. select including elements marked deleted, check most recent timestamp, to determine correct PID assignment)
 - **Add timestamp ($t-\Delta t$)** to any select statement in query
 - **Apply unique sort** to any table touched upon in query prior to query to ensure unique sort (see **R5**)

R4: Query Uniqueness

- Normalizing queries to detect identical queries
 - WHERE clause sorted
 - Calculate query string hash
 - Identify semantically identical queries
 - → non-identical queries: columns in different order



R4: Query Uniqueness





- Adapt query to history table

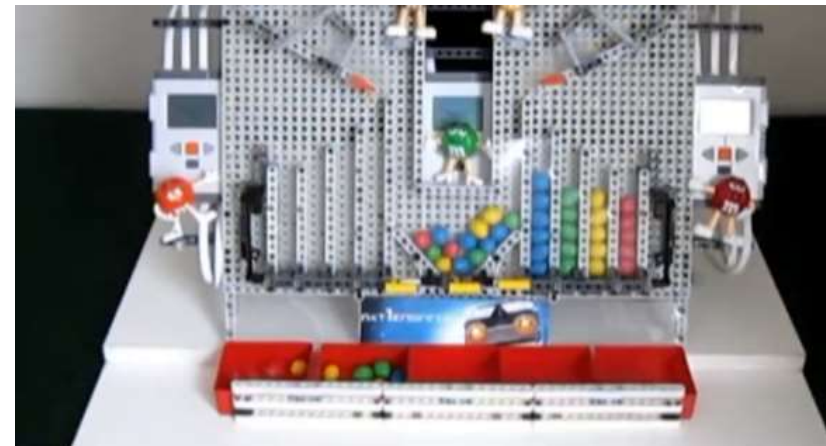
```

SELECT results.track_id, results.artist, results.release
FROM MSD AS results JOIN (
  SELECT track_id, max(timestamp) AS latestTimestamp
  FROM MSD
  WHERE timestamp <= (SELECT @queryExecutionTimestamp)
  AND (track_id NOT IN
    (SELECT track_id FROM MSD AS deletedRecords
     WHERE deletedRecords.status_mark = 'deleted'
     AND (deletedRecords.timestamp < @queryExecutionTimestamp))
  )
  GROUP BY track_id
) AS version ON results.track_id = version.track_id AND results.timestamp = version.latestTimestamp

WHERE
  results.tags = 'classic' AND results.duration > 120
ORDER BY results.track_id;
  
```


R5: Stable Sorting

- **Ensure that the sorting of the records in the data set is unambiguous and reproducible**
- The sequence of the results in the result set may not be fixed, but data processing results may depend on sequence
 - Many databases are set based
 - The storage system may use non-deterministic features
- If this needs to be addressed, apply default sort (on id)  prior to any user-defined sort
- Optional! 




<http://www.geek.com/>

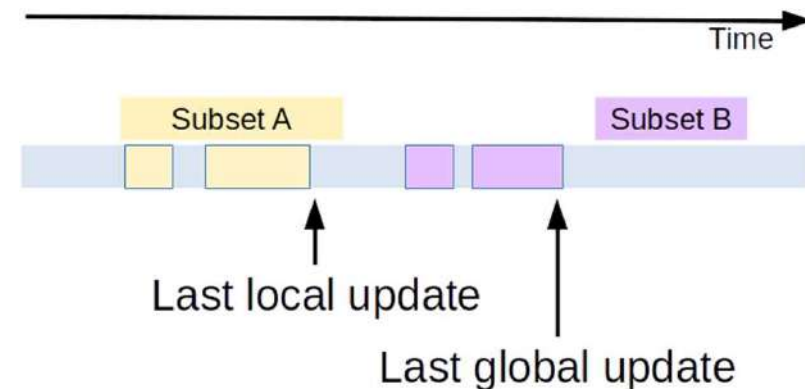
R6: Result Set Verification

- **Compute fixity information (also referred to as checksum or hash key) of the query result set to enable verification of the correctness of a result upon re-execution.**
- **Correctness:** 
 - No record has changed within a data subset
 - All records which have been in the original data set are also in the re-generated data set
- **Compute a hash key**
 - Allows to compare the completeness of results
 - For extremely large result sets: potentially limit hash input data, e.g. only row headers + record id's



R7: Query Timestamping

- Assign a timestamp to the query based on the last update to the entire database (or the last update to the selection of data affected by the query or the query execution time).
- Allows to map the execution of a query to a state of the database
 - Execution time: default solution, simple, potentially privacy concerns?
 - Last global update: simple, **recommended**
 - Last update to affected subset: complex to implement
- All equivalent in functionality!
(transparent to user) 




R8: Query PID

- **Assign a new PID to the query if either the query is new or if the result set returned from an earlier identical query is different due to changes in the data. Otherwise, return the existing PID of the earlier query to the user.**
- **Existing PID:** Identical query (semantics) with identical result set, i.e. no change to any element touched upon by query since first processing of the query
- **New PID:** whenever query semantics is not absolutely identical
(irrespective of result set being potentially identical!)

R8: Query PID

- Note:
 - Identical result set alone does not mean that the query semantics is identical
 - Will assign different PIDs to capture query semantics
 - Need to normalize query to allow comparison
- Process:
 - Re-write query to adapt to versioning system, stable sorting, ...
 - Determine query hash
 - Execute user query and determine result set hash
 - Check query store for queries with identical query hash
 - If found, check for identical result set hash
- 2 PIDs: (compare e.g. paper in journal)
 - precise subset of (static) data, as an excerpt of
 - a larger, dynamically evolving data stream

R9: Store the Query

- **Store query and metadata (e.g. PID, original and normalised query, query and result set checksum, timestamp, superset PID, data set description, and other) in the query store.**
 - Query store is central infrastructure
 - Stores query details for long term
 - Provides information even when the data should be gone
 - Responsible for re-execution
 - Holds data for landing pages
 - Stores sensitive information
- Not necessarily ALL queries (staging area) 



R10: Create Citation Texts

- **Generate citation texts in the format prevalent in the designated community for lowering the barrier for citing and sharing the data. Include the PID in the citation text snippet.**
- Researchers are “lazy”/efficient
 - Support citing by allow them to copy and paste citations for data
 - Citations contain text including PIDs and timestamps
 - Adapted for each community
- **2 PIDs!**
 - Superset: the “database” and it’s holder (repository, data center)
 - Subset: based on the query
 - Accumulate credits for subset and (dynamic) data collection/holder

Suggested citation text:

Stefan Proell (2015) "Austria Facts" created at 2015-10-07 10:51:55.0, PID [ark:12345/qmZi2wO2vy]. Subset of CIA: "The CIA WorldFactbook", PID [ark:12345/cLfH9FjxnA]

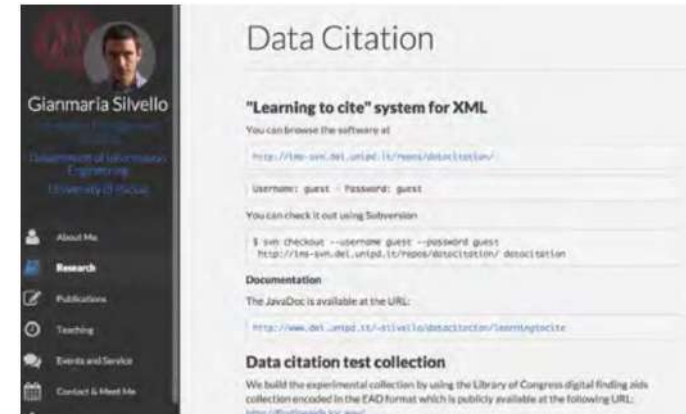


R10: Automated Citation Texts

- Can be created automatically
 - relatively simple for relational
 - more complex for hierarchical/XML

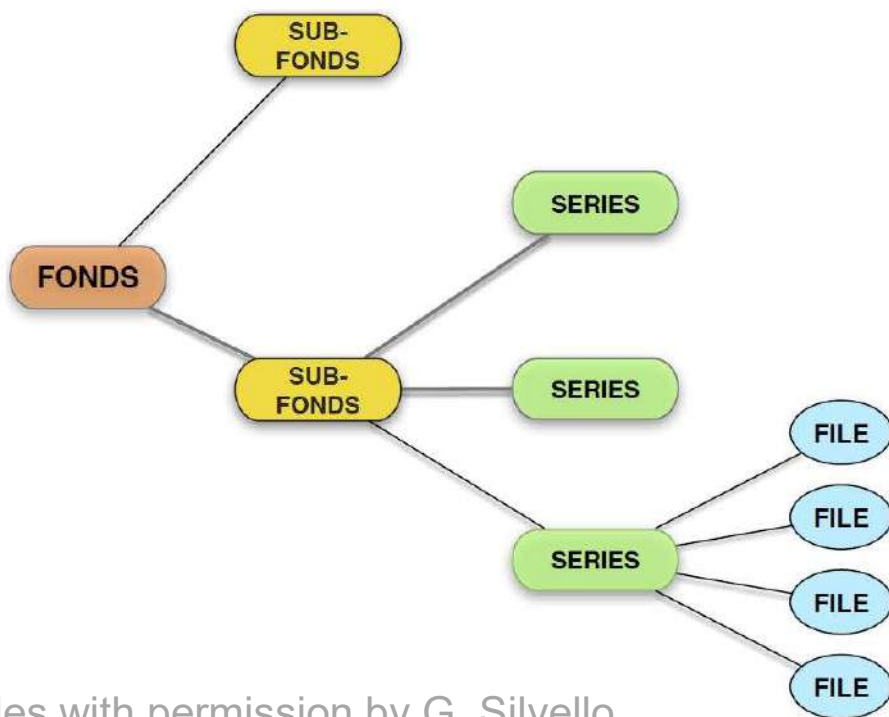
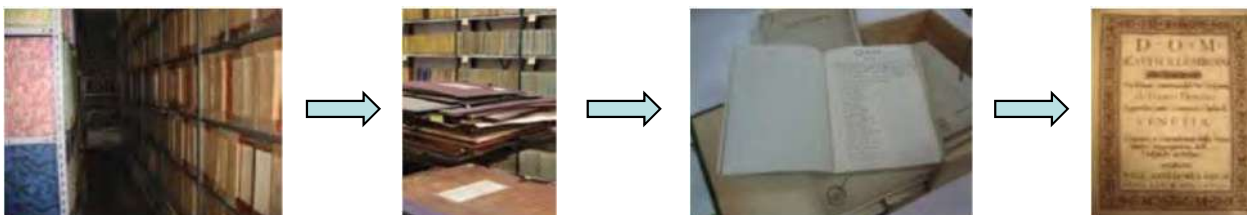
- Learning to Cite:

- Gianmaria Silvello. Learning to Cite Framework: How to Automatically Construct Citations for Hierarchical Data. Journal of the Association for Information Science and Technology (JASIST), Volume 68 issue 6, pp. 1505-1524, June 2017.
- <http://www.dei.unipd.it/~silvello/datacitation>



R10: Automated Citation Texts

- EAD: Encoded Archival Description



```

<ead>
  <eadheader>
    [...]
  </eadheader>
  <archdesc level="fonds">
    [...]
    <did>[...]</did>
    <dsc level="fonds">
      [...]
      <c01 level="sub-fonds">
        [...]
      </c01>
      <c01 level="sub-fonds">
        [...]
        <c02 level="series">
          [...]
        </c02>
        <c02 level="series">
          [...]
        </c02>
        <c02 level="series">
          [...]
          <c03 level="file">
            [...]
          </c03>
          <c03 level="file">
            [...]
          </c03>
          <c03 level="file">
            [...]
          </c03>
          <c03 level="file">
            [...]
          </c03>
        </c02>
      </c01>
    </dsc>
  </archdesc>
</ead>
  
```



R10: Automated Citation Texts

- A human-readable citation:

Correspondence, 1951-1956,

"The Elements of Legal Theory" (unpublished). Books, box 135. Part II:

Writings (1905-1984), box 129-152. Huntington Cairns Papers.

Manuscript Division, Library of Congress.

<http://hdl.loc.gov/loc.mss/eadmss.ms001024>

R10: Automated Citation Texts

- A human-readable citation:

Citable unit

Correspondence, 1951-1956

Contextual Information (from ancestors of the citable unit)

"The Elements of Legal Theory" (unpublished). Books, box 135. Part II:
Writings (1905-1984), box 129-152. Huntington Cairns Papers.
Manuscript Division, Library of Congress.

<http://hdl.loc.gov/loc.mss/eadmss.ms001024>

(Persistent) Unique identifier of the EAD file



R10: Automated Citation Texts

- A machine-readable citation:
 - Conjunction of XML paths

```
/ead/eadheader/eadid && /ead/eadheader/filedesc/publicationstmt/publisher && /ead/  
archdesc/did/unittitle && /ead/archdesc/dsc/c01[10]/did/unittitle && /ead/archdesc/  
dsc/c01[10]/did/unittitle/unitdate && /ead/archdesc/dsc/c01[10]/did/container/@type  
&& /ead/archdesc/dsc/c01[10]/did/container && /ead/archdesc/dsc/c01[10]/c02/did/  
container/@type && /ead/archdesc/dsc/c01[10]/c02/did/container && /ead/archdesc/dsc/  
c01[10]/c02/did/unittitle && /ead/archdesc/dsc/c01[10]/c02/c03[4]/did/unittitle  
&& /ead/archdesc/dsc/c01[10]/c02/c03[4]/did/container/@type && /ead/archdesc/dsc/  
c01[10]/c02/c03[4]/did/container && /ead/archdesc/dsc/c01[10]/c02/c03[4]/c04[2]/  
did/unittitle && /ead/archdesc/dsc/c01[10]/c02/c03[4]/c04[2]/c05[1]/did/unittitle
```

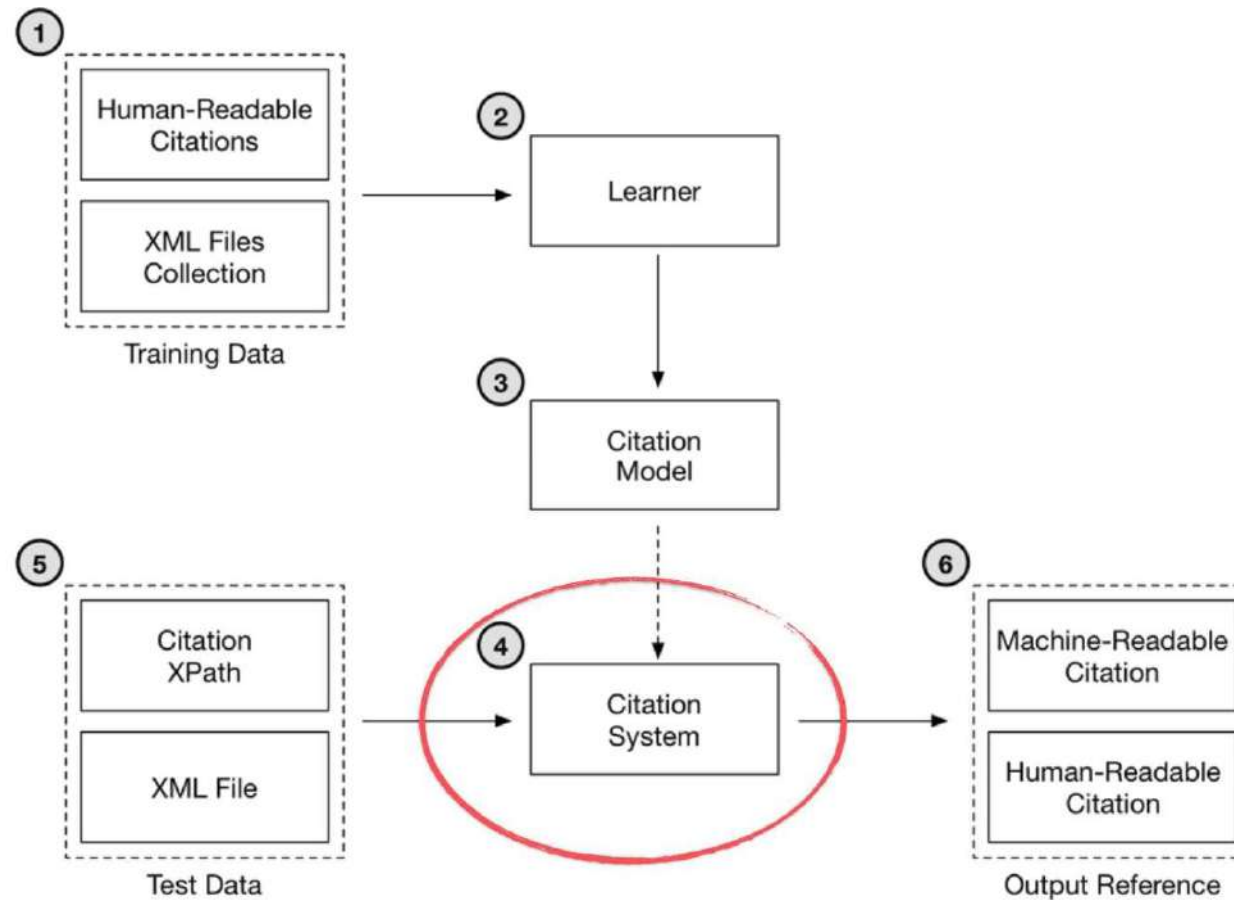
R10: Automated Citation Texts

- Mapping machine-readable to human-readable:

Human-Readable Citation	Machine-Readable Citation
http://hdl.loc.gov/loc.mss/eadmss.ms001024 ←	/ead/eadheader/eadid
Manuscript Division, Library of Congress ←	/ead/eadheader/filedesc/publicationstmt/publisher
Huntington Cairns Papers ←	/ead/archdesc/did/unittitle
Part II: Writings ←	/ead/archdesc/dsc/c01[10]/did/unittitle
1905-1984 ←	/ead/archdesc/dsc/c01[10]/did/unittitle/unitdate
box ←	/ead/archdesc/dsc/c01[10]/did/container/@type
129-152 ←	/ead/archdesc/dsc/c01[10]/did/container
By Cairns ←	/ead/archdesc/dsc/c01[10]/c02[1]/did/unittitle
box ←	/ead/archdesc/dsc/c01[10]/c02[1]/did/container/@type
129 ←	/ead/archdesc/dsc/c01[10]/c02[1]/did/container/
Books ←	/ead/archdesc/dsc/c01[10]/c02[1]/c03[4]/did/unittitle
box ←	/ead/archdesc/dsc/c01[10]/c02[1]/c03[4]/did/container/@type
135 ←	/ead/archdesc/dsc/c01[10]/c02[1]/c03[4]/did/container
"The Elements of Legal Theory" (unpublished) ←	/ead/archdesc/dsc/c01[10]/c02[1]/c03[4]/c04[2]/did/unittitle
Correspondence, 1951-1956 ←	/ead/archdesc/dsc/c01[10]/c02[1]/c03[4]/c04[2]/c05[1]/did/unittitle

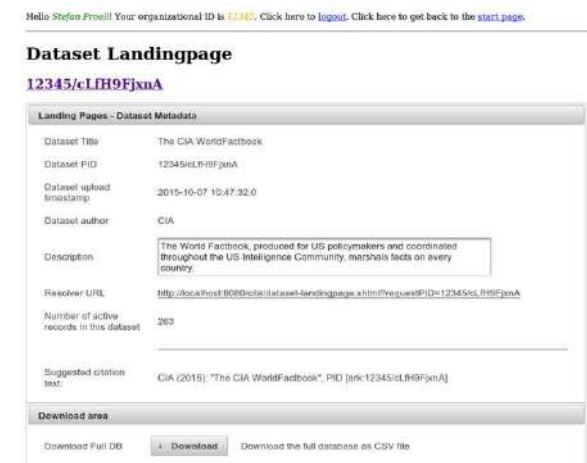
R10: Automated Citation Texts

- Learning citation models



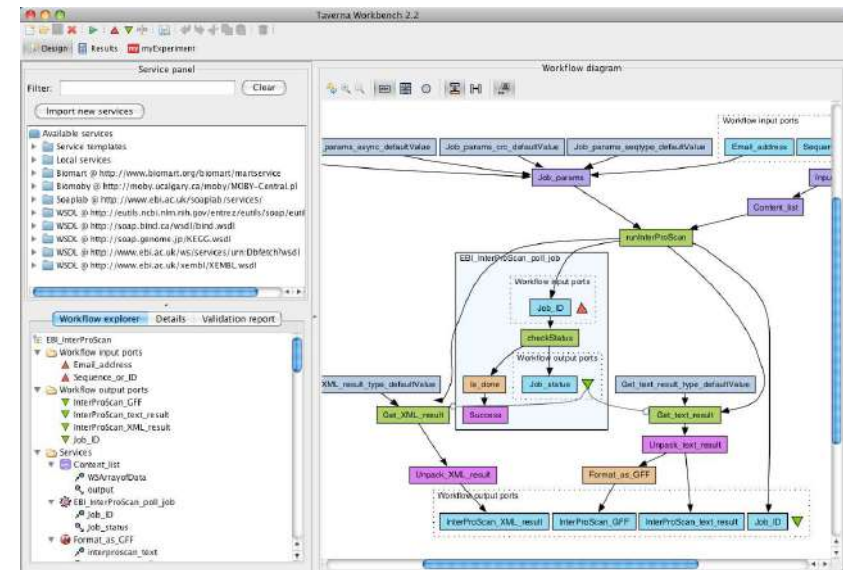
R11: Landing Page

- **Make the PIDs resolve to a human readable landing page that provides the data (via query re-execution) and metadata, including a link to the superset (PID of the data source) and citation text snippet.**
 - Data sets and subsets uniquely identifiable by their PID, which resolves to a human readable landing page.
 - Landing page reachable by a unique URL, presented in a Web browser
 - Not all information needs to be provided on landing page (e.g. query strings frequently not relevant / potential security threat)




R12: Machine Actionability

- **Provide an API / machine actionable landing page to access metadata and data via query re-execution.**
 - Experiments are increasingly automated
 - Machines most likely to consume data citations
 - Allows machines to resolve PIDs, access metadata and data
 - Note: does NOT imply full / automatic access to data!
 - Authentication
 - Load analysis
 - Handshake, content negotiation, ...
 - Allows automatic meta-studies, monitoring, ...



R13: Technology Migration

- **When data is migrated to a new representation (e.g. new database system, a new schema or a completely different technology), migrate also the queries and associated fixity information.**
 - Technology evolves and data may be moved to a new technology stack
 - Query languages change
- Migration required 
 - Migrate data and the queries (both are with the data center!)
 - Adapt versioning, re-compute query hash-keys
 - Maybe decide to keep “original” queries in the provenance trace
- Note: such data migrations constitute major projects, usually happen rarely – require all APIs to be adapted, ...

R13: Technology Migration

- **Consider e.g. Schema Modification Operators (SMOs)**
 - CREATE TABLE R
 - DROP TABLE R
 - RENAME TABLE R
 - COPY TABLE R INTO S
 - PARTITION TABLE R INTO S with *cond*, T with *!cond*
 - DECOMPOSE TABLE R INTO S(A,B) T(A,C)
 - JOIN TABLE R,S INTO T WHERE *cond*
 - ADD COLUMN C [as const | func(A)] INTO R
 - DROP COLUMN C FROM R
 - RENAME COLUMN A IN R TO B
- How will the queries need to be re-written to address them?



R14: Migration Verification

- **Verify successful data and query migration, ensuring that queries can be re-executed correctly.**
- Sanity check: After migration is done, verify that the data can still be retrieved correctly
- Use query and result set hashes in the query store to verify results
- If hash function is incompatible/cannot be computed on new system as hash input data sequence cannot be obtained, pairwise comparison of subset elements
 - May constitute new PID / data subset in this case, as subsequent processes will not be able to use it as input if result set presentation has changed, breaks processes

RDA Recommendations - Summary

- Building blocks of supporting dynamic data citation:
 - Uniquely identifiable data records
 - Versioned data, marking changes as insertion/deletion
 - Time stamps of data insertion / deletions
 - “Query language” for constructing subsets
- Add modules:
 - Persistent query store: queries and the timestamp (either: <when issued> or <of last change to data>)
 - Query rewriting module
 - PID assignment for queries that enables access
- Stable across data source migrations (e.g. diff. DBMS), scalable, machine-actionable

RDA Recommendations - Summary

▪ *Benefits*

- Allows **identifying, retrieving and citing the precise data subset** with minimal storage overhead by only storing the versioned data and the queries used for extracting it
- Allows retrieving the data both **as it existed** at a given point in time as well as the **current view** on it, by re-executing the same query with the stored or current timestamp
- It allows to cite even an **empty set!**
- The query stored for identifying data subsets provides valuable **provenance data**
- Query store collects **information on data usage**, offering a basis for data management decisions
- **Metadata** such as checksums support the verification of the correctness and **authenticity** of data sets retrieved
- The same principles work for **all types of data**

- ***Some considerations and questions***

- May data be deleted?

Yes, of course, given appropriate policies. Queries may then not be re-executable against the original timestamp anymore

- Does the system need to store every query?

No, only data sets that should be persisted for citation and later re-use need to be stored.

- Can I obtain only the most recent data set?

Queries can be re-executed with the original timestamp or with the current timestamp or any other timestamp desired.

- Which PID system should be used?


Any PID system can, in principle, be applied according to the institutional policy.

Outline

-
- Why should we want to cite data?
 - What identifier system should I use?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - Summary
-

Large Number of Adoptions

■ Standards / Reference Guidelines / Specifications:

- Joint Declaration of Data Citation Principles: 
Principle 7: Specificity and Verifiability (<https://www.force11.org/datacitation>)
- ESIP:Data Citation Guidelines for Earth Science Data Vers. 2 (P14)
- ISO 690, Information and documentation - Guidelines for bibliographic references and citations to information resources (P13)
- EC ICT TS5 Technical Specification (pending) (P12)
- DataCite Considerations (P8)

■ Reference Implementations

- MySQL/Postgres (P5, P6)
- CSV files: MySQL, Git (P5, P6, P8, Webinar)
- XML (P5)
- CKAN Data Repository (P13)
- RDF/SPARQL (P17)

Large Number of Adoptions

- **Pilot implementations, Use cases**
 - DEXHELPP: Social Security Records (P6)
 - NERC: ARGO Global Array (P6)
 - LNEC: River dam monitoring (P5)
 - CLARIN: Linguistic resources, XML (P5)
 - MSD: Million Song Database (P5)
 - many further individual ones discussed ...

Large Number of Adoptions

■ Adoptions deployed

- CBMI: Center for Biomedical Informatics, WUSTL (P8, Webinar)
- VMC: Vermont Monitoring Cooperative (P8, Webinar)
- CCCA: Climate Change Center Austria (P10/P11/P12, Webinar)
- EODC: Earth Observation Data Center (P14, Webinar)
- VAMDC: Virtual Atomic and Molecular Data Center (P8/P10/P12, Webinar)
- Ocean Networks Canada (P12, Webinar)

■ In progress

- NICT Smart Data Platform (P10/P14)
- Dendro System (P13)
- Deep Carbon Observatory (P12)



WGDC Webinar Series

- <https://www.rd-alliance.org/group/data-citation-wg/webconference/webconference-data-citation-wg.html>
 - Implementation of the RDA Data Citation Recommendations by **Ocean Networks Canada (ONC)**
 - Implementation of the RDA Data Citation Recommendations the **Earth Observation Data Center (EODC) for the openEO platform**
 - **Automatically generating citation text from queries for RDBMS and XML data sources**
 - Implementing of the RDA Data Citation Recommendations by the **Climate Change Centre Austria (CCCA) for a repository of NetCDF files**
 - Implementing the RDA Data Citation Recommendations for **Long-Tail Research Data / CSV files**
 - Implementing the RDA Data Citation Recommendations in the **Distributed Infrastructure of the Virtual and Atomic Molecular Data Center (VAMDC)**
 - Implementation of Dynamic Data Citation at the **Vermont Monitoring Cooperative**
 - Adoption of the RDA Data Citation of Evolving Data Recommendation to **Electronic Health Records**

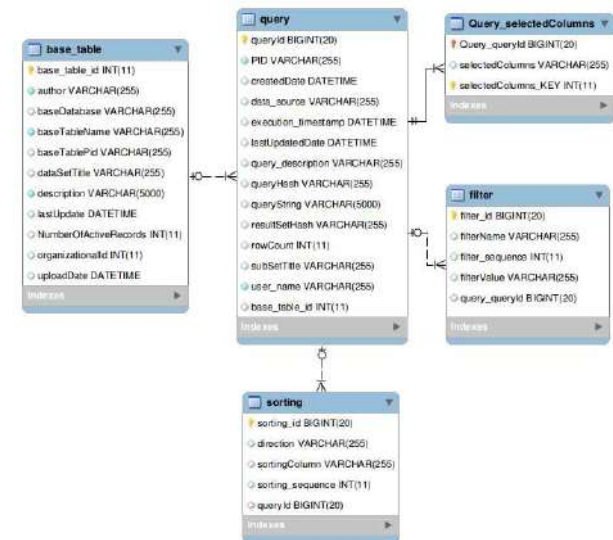




**Reference Implementation for
CSV Data (and SQL)
Stefan Pröll, SBA
Christoph Meixner, TU Wien**

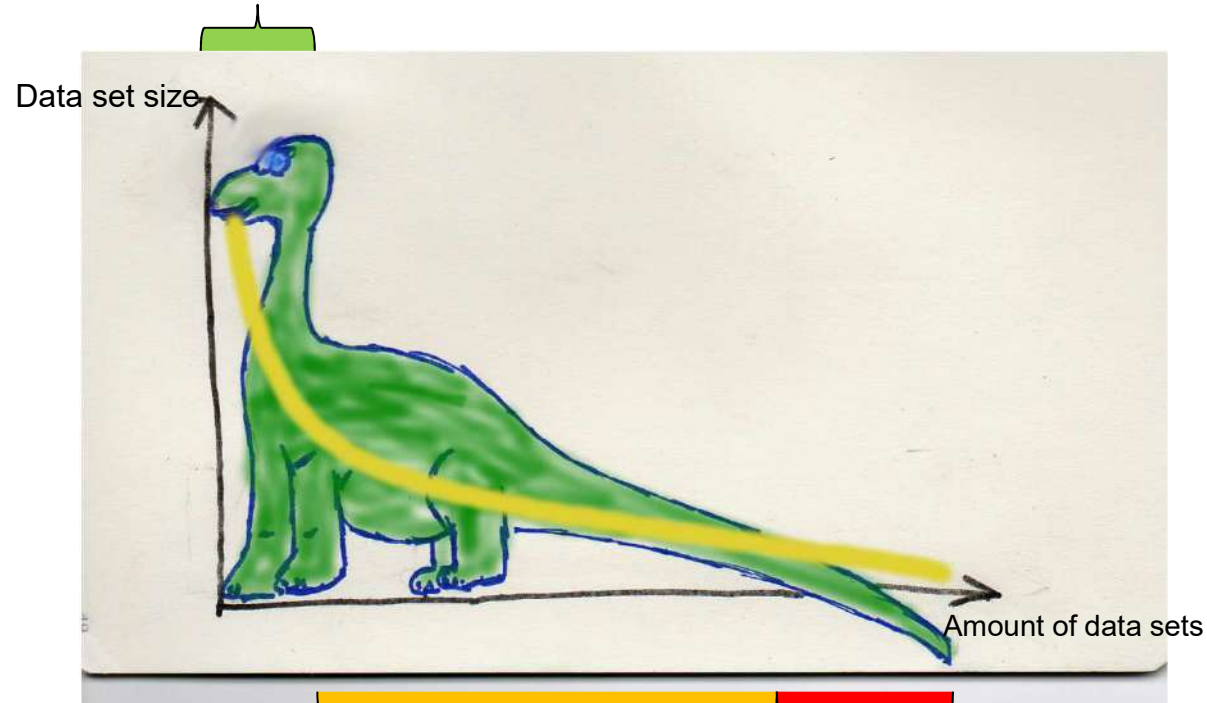
research data sharing without barriers
rd-alliance.org

- RDA recommendations implemented in data infrastructures
- Required adaptations
 - Introduce versioning, if not already in place
 - Capture sub-setting process (queries)
 - Implement dedicated query store to store queries
 - A bit of additional functionality (query re-writing, hash functions, ...)
- Done! ?
 - “Big data”, database driven
 - Well-defined interfaces
 - Trained experts available
 - “Complex, only for professional research infrastructures” ?



Long Tail Research Data

Big data,
well organized,
often used and cited



Less well organized, “Dark data”
non-standardised
no dedicated infrastructure

Prototype Implementations

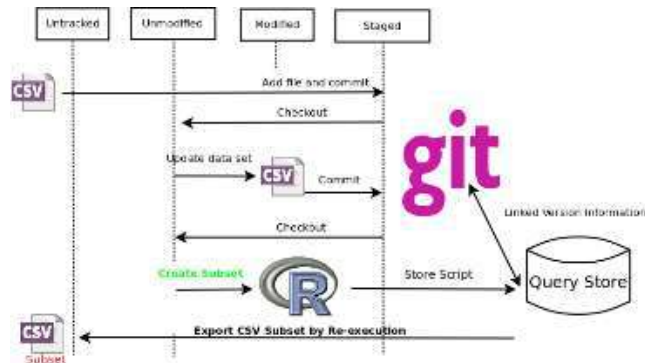
- Solution for small-scale data
 - CSV files, no “expensive” infrastructure, low overhead

2 Reference implementations :

- **Git** based Prototypes: widely used versioning system
 - A) Using separate folders
 - B) Using branches
- **MySQL** based Prototype:
 - C) Migrates CSV data into relational database
- Data backend responsible for versioning data sets
- Subsets are created with scripts or queries via API or Web Interface
- Transparent to user: always CSV

Git Implementation 1

- Upload CSV files to Git repository (versioning)
- Subsets created via scripting language (e.g. R)
 - Select rows/columns, sort, returns CSV + metadata file
 - Metadata file with script parameters stored in Git
 - (Scripts stored in Git as well)
- PID assigned to metadata file
 - Use Git to retrieve proper data set version and re-execute script on retrieved file



```
diff --git a/suppList.csv b/suppList.csv
index 7e2799f..95b0225 100644
--- a/suppList.csv
+++ b/suppList.csv
@@ -1,21 +1,21 @@
#suppName,sex,weight,height,eyeColor,smoke
0,Cat,20.0,0.28438289328562742,20340,gray,white,non-smoker
1,cat,sm,0.28438289328562742,20340,gray,white,non-smoker
2,cat,sm,0.43251899328562742,20340,gray,white,non-smoker
3,cat,sm,0.494899328562742,20340,gray,white,non-smoker
4,cat,sm,0.50274562742,20340,gray,white,non-smoker
5,cat,sm,0.50274562742,20340,gray,white,non-smoker
6,cat,sm,0.50274562742,20340,gray,white,non-smoker
7,cat,sm,0.50274562742,20340,gray,white,non-smoker
8,cat,sm,0.50274562742,20340,gray,white,non-smoker
9,cat,sm,0.50274562742,20340,gray,white,non-smoker
10,cat,sm,0.50274562742,20340,gray,white,non-smoker
11,cat,sm,0.50274562742,20340,gray,white,non-smoker
12,cat,sm,0.50274562742,20340,gray,white,non-smoker
13,cat,sm,0.50274562742,20340,gray,white,non-smoker
14,cat,sm,0.50274562742,20340,gray,white,non-smoker
15,cat,sm,0.50274562742,20340,gray,white,non-smoker
16,cat,sm,0.50274562742,20340,gray,white,non-smoker
17,cat,sm,0.50274562742,20340,gray,white,non-smoker
18,cat,sm,0.50274562742,20340,gray,white,non-smoker
19,cat,sm,0.50274562742,20340,gray,white,non-smoker
```

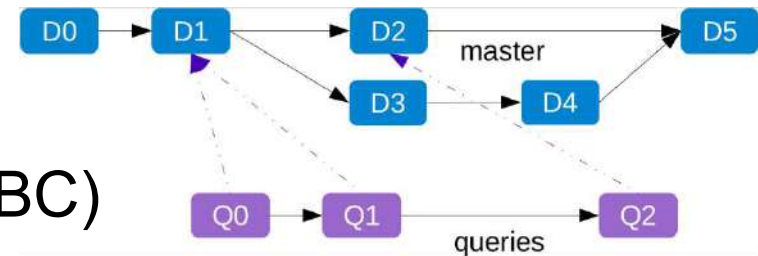
```
# PID=1234/abcd5fgh
# Repository.Path=/media/Data/Git-Repository
# Execution.Time=2015-09-30:11:07:09
# Subset.Tool=scripting front-end version 3.0.2 (2015-08-14)
# Subset.Tool.Path=/usr/bin/Rscript
# Input.Script.Path=/supercomputing/top5-script.v
# Input.Script.Hash=baf5d...d7881:supercomputing/top5-script.r
# Dataset.Path=/supercomputing/supercomputer.csv
# Dataset.Commit.Hash=caed...d788c:supercomputer.csv
# Output.Path=/tmp/supercomputer-top5.csv

# Original execution:
# /usr/bin/Rscript /supercomputing/top5-script.v \
# /media/Data/Git-Repository/supercomputing/supercomputer.csv \
# /tmp/supercomputer-top5.csv

# Recommended re-execution
# Retrieve script
git --git-dir=/media/Data/Git-Repository/.git/ \
show baf5d...d7881:supercomputing/top5-script.r
> /tmp/reproduced-datasets/top5-script.r
# Retrieve data set
git --git-dir=/media/Data/Git-Repository/.git/ \
show d78ed...b9792:supercomputing/supercomputer.csv \
> /tmp/reproduced-datasets/supercomputer.csv
# Re-execute
/usr/bin/Rscript /supercomputing/top5-script.r \
/tmp/reproduced-datasets/supercomputer.csv \
/tmp/reproduced-datasets/supercomputer-top5.csv
```

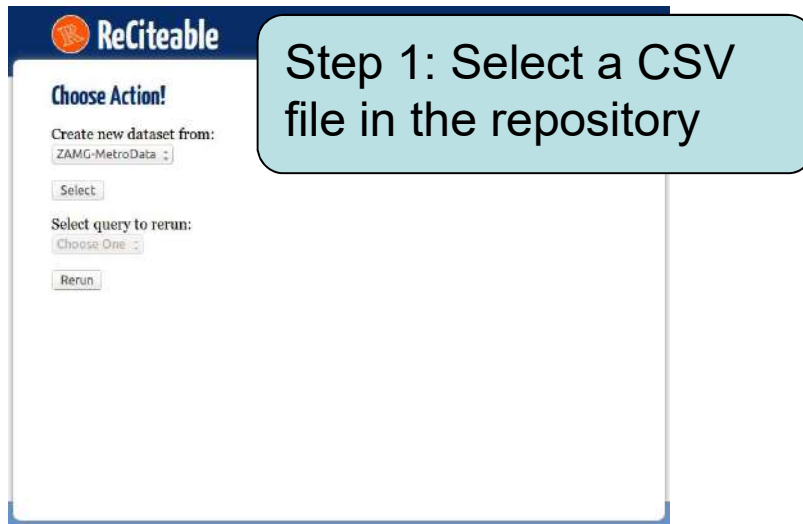
Git Implementation 2

- Addresses issues
 - common commit history, branching data
- Using Git branching model:
Orphaned branches for queries and data
 - Keeps commit history clean
 - Allows merging of data files
- Web interface for queries (CSV2JDBC)
- Use commit hash for identification
 - Assigned PID hashed with SHA1
 - Use hash of PID as filename (ensure permissible characters)



Git-Based Reference Implementation

- Prototype: <https://github.com/Mercynary/recitable>



Step 1: Select a CSV file in the repository

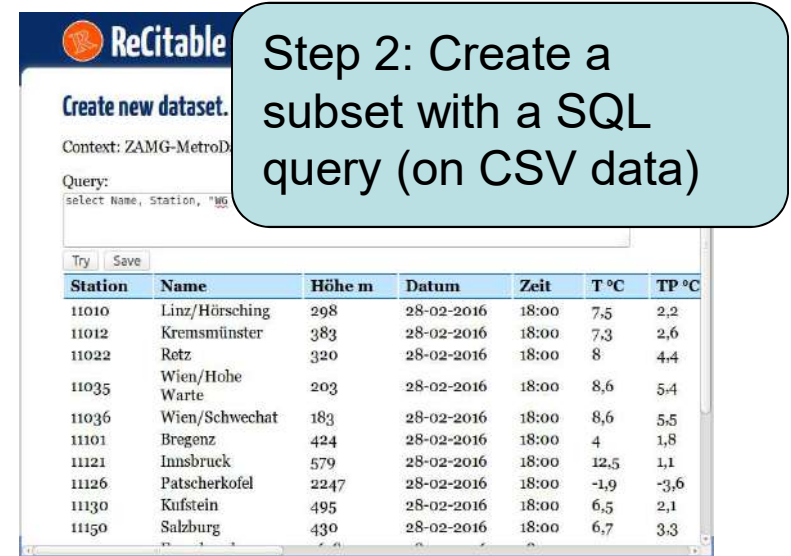
Choose Action!

Create new dataset from:
ZAMG-MetroData

Select

Select query to rerun:
Choose One

Rerun



Step 2: Create a subset with a SQL query (on CSV data)

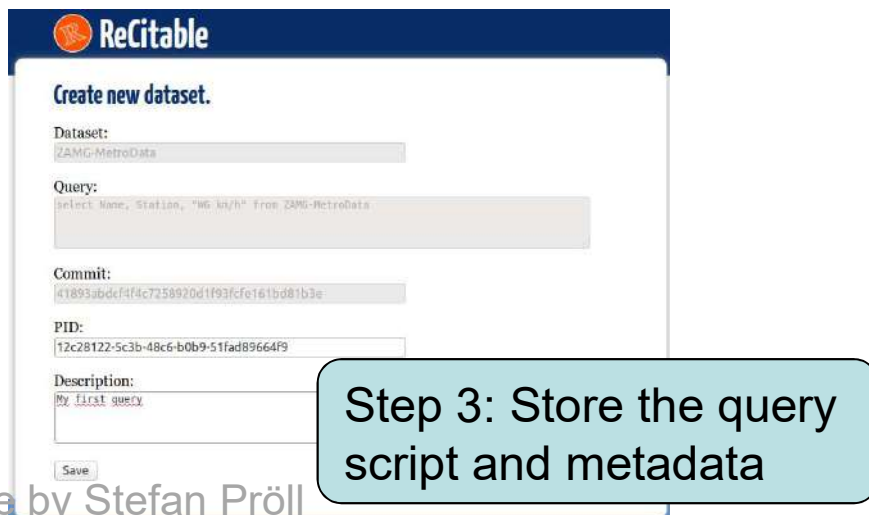
Create new dataset.

Context: ZAMG-MetroData

Query:
Select Name, Station, "HG

Try Save

Station	Name	Höhe m	Datum	Zeit	T °C	TP °C
11010	Linz/Hörsching	298	28-02-2016	18:00	7,5	2,2
11012	Kremsmünster	383	28-02-2016	18:00	7,3	2,6
11022	Retz	320	28-02-2016	18:00	8	4,4
11035	Wien/Hohe Warte	203	28-02-2016	18:00	8,6	5,4
11036	Wien/Schwechat	183	28-02-2016	18:00	8,6	5,5
11101	Bregenz	424	28-02-2016	18:00	4	1,8
11121	Innsbruck	579	28-02-2016	18:00	12,5	1,1
11126	Patscherkofel	2247	28-02-2016	18:00	-1,9	-3,6
11130	Kufstein	495	28-02-2016	18:00	6,5	2,1
11150	Salzburg	430	28-02-2016	18:00	6,7	3,3



Step 3: Store the query script and metadata

Create new dataset.

Dataset:
ZAMG-MetroData

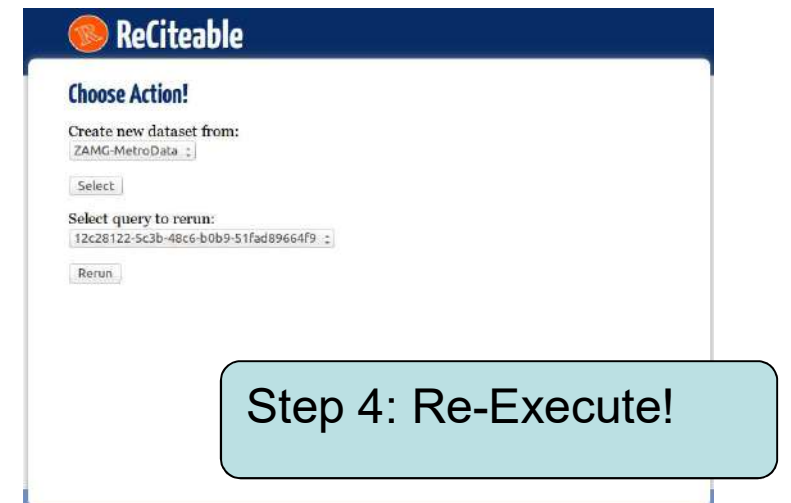
Query:
Select Name, Station, "HG m/h" From ZAMG-MetroData

Commit:
81893abdef1f4c7258920d1f93cfe161bd81b3e

PID:
12c28122-5c3b-48c6-b0b9-51fad89664f9

Description:
My first query

Save



Step 4: Re-Execute!

Choose Action!

Create new dataset from:
ZAMG-MetroData

Select

Select query to rerun:
12c28122-5c3b-48c6-b0b9-51fad89664f9

Rerun

MySQL Prototype

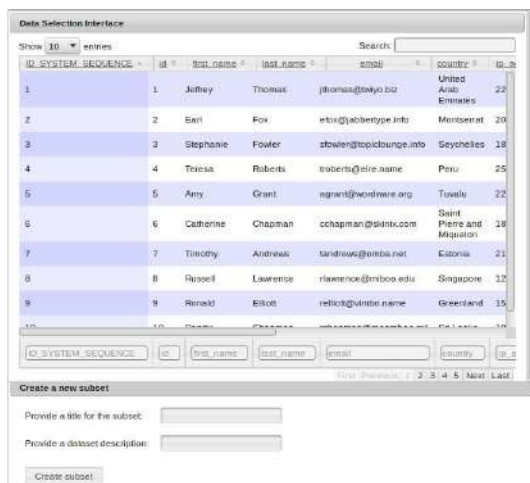
- Data upload
 - User uploads a CSV file into the system
- Data migration from CSV file into RDBMS
 - Generate table structure
 - Add metadata columns (versioning)
 - Add indices (performance)
- Dynamic data
 - Insert, update and delete records
 - Events are recorded with a timestamp
- Subset creation
 - User selects columns, filters and sorts records in web interface
 - System traces the selection process
 - Exports CSV





MySQL-Based Reference Implementation

- Source at Github:
 - <https://github.com/datascience/RDA-WGDC-CSV-Data-Citation-Prototype>
- Videos:
 - Login: <https://youtu.be/EnralwbQfM0>
 - Upload: <https://youtu.be/xJruifX9E2U>
 - Subset: <https://www.youtube.com/watch?v=it4sC5vYiZQ>
 - Resolver: <https://youtu.be/FHsvjsUMiiY>
 - Update: <https://youtu.be/cMZ0xoZHUyI>



CSV Reference Implementations

- Stefan Pröll, Christoph Meixner, Andreas Rauber
Precise Data Identification Services for Long Tail Research Data.
Proceedings of the intl. Conference on Preservation of Digital Objects
(iPRES2016), Oct. 3-6 2016, Bern, Switzerland.

- Source at Github:
<https://github.com/datascience/RDA-WGDC-CSV-Data-Citation-Prototype>

Videos:

- Login: <https://youtu.be/EnralwbQfM0>
- Upload: <https://youtu.be/xJruifX9E2U>
- Subset: <https://www.youtube.com/watch?v=it4sC5vYiZQ>
- Resolver: <https://youtu.be/FHsvjsUMiiY>
- Update: <https://youtu.be/cMZ0xoZHUyl>





WG Data Citation Pilot CBMI @ WUSTL

**Cynthia Hudson Vitale, Leslie McIntosh,
Snehil Gupta
Washington University in St. Luis**

research data sharing without barriers
rd-alliance.org

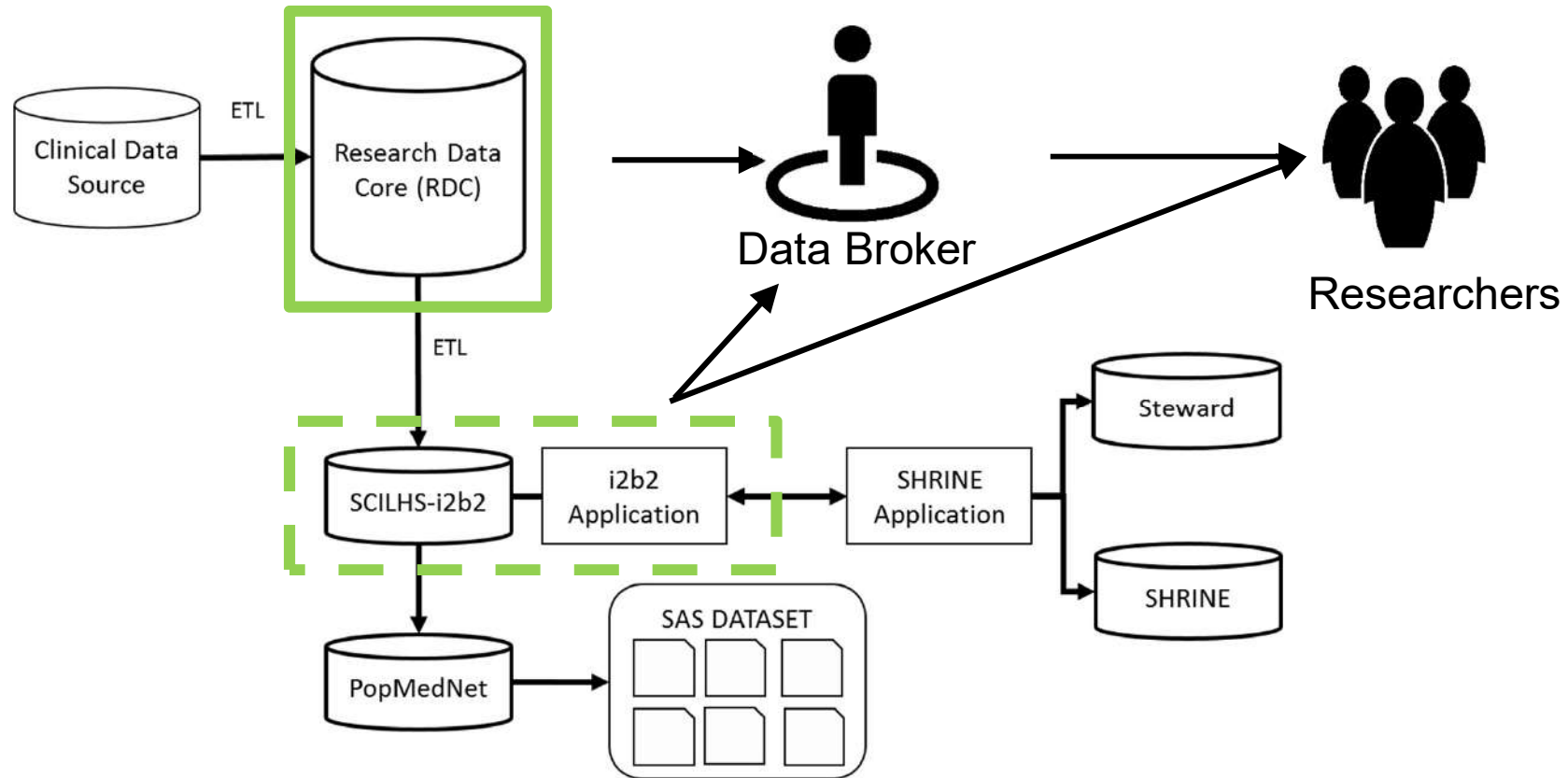


Biomedical Adoption Project Goals

- Implement RDA Data Citation WG recommendation to local Washington U i2b2
- Engage other i2b2 community adoptees
- Contribute source code back to i2b2 community

- Repository
https://github.com/CBMIWU/Research_Reproducibility
- Slides
<http://bit.ly/2cnWorU>
- Bibliography
https://www.zotero.org/groups/biomedical_informatics_resrepro

RDA-MacArthur Grant Focus



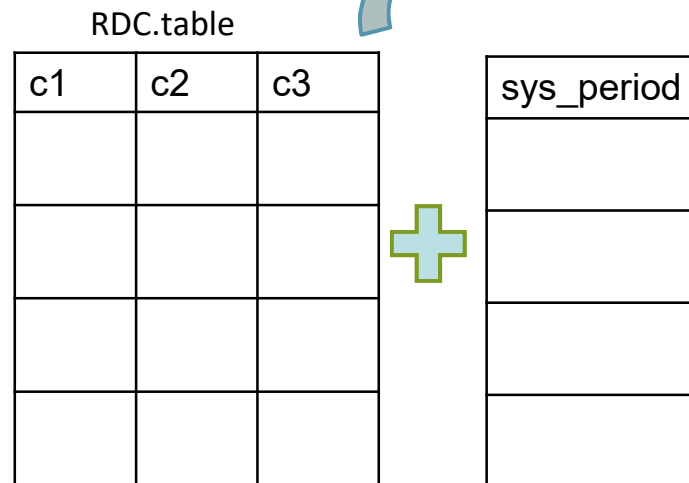
R1 and R2 Implementation

1



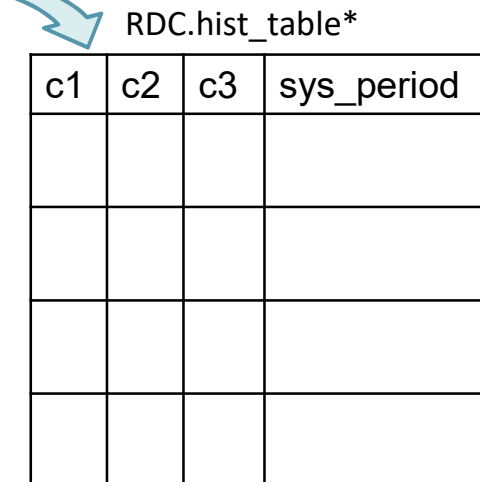
PostgreSQL Extension
"temporal_tables"

2



triggers

3



*stores history of data changes



Return on Investment (ROI) - Estimated

- 20 hours to complete 1 study
- \$150/hr (unsubsidized)
- \$3000 per study
- 115 research studies per year
- **14 replication studies**



CBMI Research Reproducibility Resources

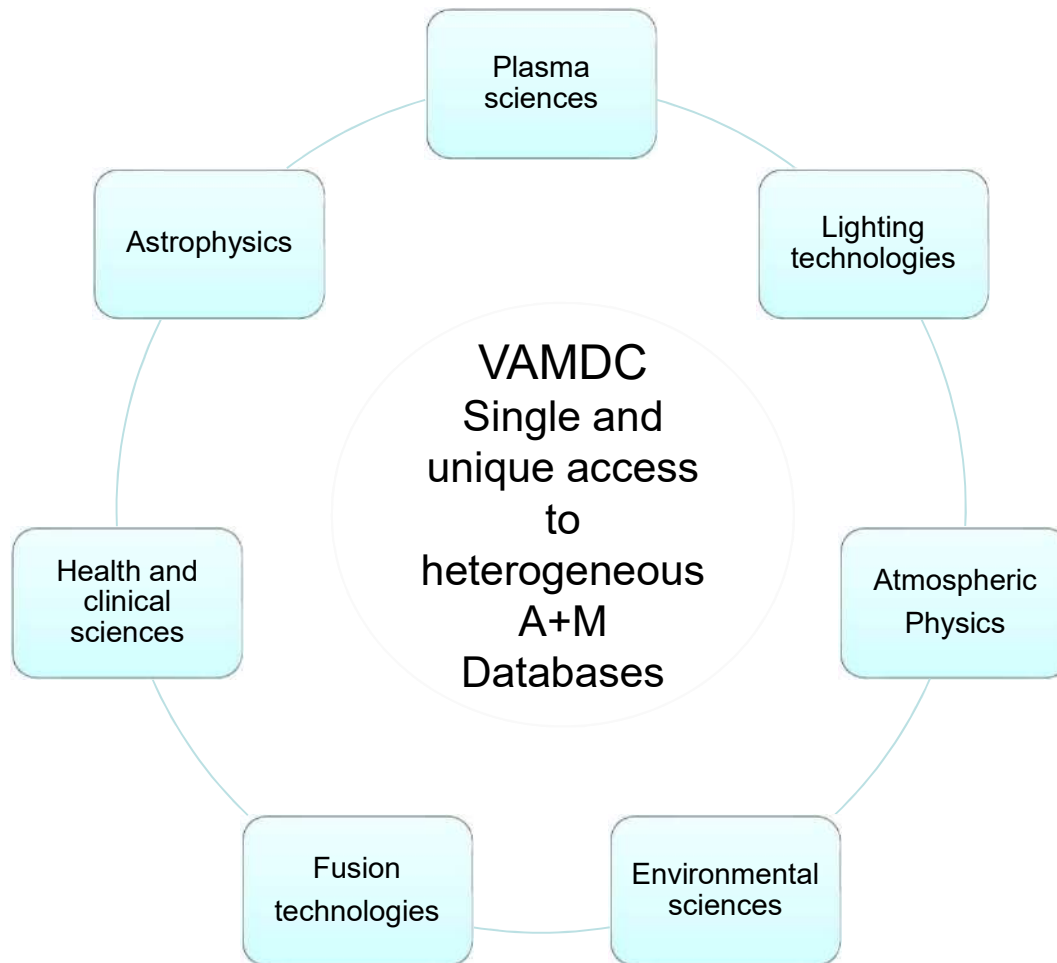
-
- Repository
https://github.com/CBMiWU/Research_Reproducibility
 - Slides
<http://bit.ly/2cnWorU>
 - Bibliography
https://www.zotero.org/groups/biomedical_informatics_resepro



**From RDA Data Citation
Recommendations to new paradigms for
citing data from VAMDC
C.M. Zwölf and VAMDC Consortium
*carlo-maria.zwolf@obspm.fr***

research data sharing without barriers
rd-alliance.org

The Virtual Atomic and Molecular Data Centre



➤ Federates 29 heterogeneous databases
<http://portal.vamdc.org/>

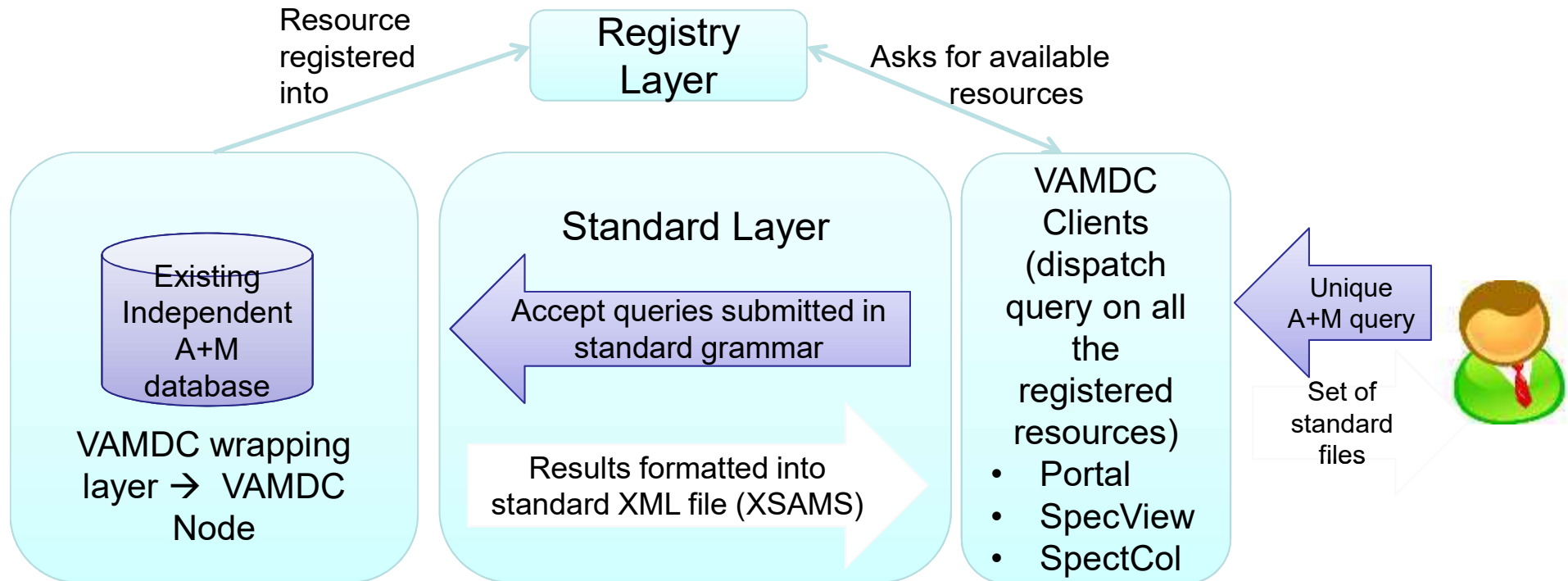
➤ The “V” of VAMDC stands for Virtual in the sense that the e-infrastructure does not contain data. The infrastructure is a wrapping for exposing in a unified way a set of heterogeneous databases.

➤ The consortium is politically organized around a Memorandum of understanding (15 international members have signed the MoU, 1 November 2014)

➤ High quality scientific data come from different Physical/Chemical Communities

➤ Provides data producers with a large dissemination platform

➤ Remove bottleneck between data-producers and wide body of users



- VAMDC is agnostic about the local data storage strategy on each node.
- Each node implements the access/query/result protocols.
- There is no central management system.
- Decisions about technical evolutions are made by consensus in Consortium.

➤ It is both technical and political challenging to implement the WG recommendations.

Let us implement the recommendation!!

Implementation will be an overlay to the standard / output layer, thus independent from any specific data-node

Tagging versions of data

Two layers
mechanisms

1 → Fine grained granularity:
Evolution of XSAMS output
standard for tracking data
modifications

2 → Coarse grained
granularity:
At each data modification to a
given data node, the version
of the Data-Node changes

With the second mechanism we know that something changed : in other words, we know that the result of an identical query may be different from one version to the other. The detail of which data changed is accessible using the first mechanisms.

Let us implement the recommendation!!

Implementation will be an overlay to the standard / output layer, thus independent from any specific data-node

Tagging versions of data

Query Store

Two layers
mechanisms

1 → **Fine grained granularity:**
Evolution of XSAMS output standard for tracking data modifications

2 → **Coarse grained granularity:**

At each data modification to a given data node, the version of the Data-Node changes

With the second mechanism we know that something changed : in other words, we know that the result of an identical query may be different from one version to the other. The detail of which data changed is accessible using the first mechanisms

Is built over the versioning of Data

Is plugged over the existing VAMDC data-extraction mechanisms

Due to the distributed VAMDC architecture, the Query Store architecture is similar to a log-service



Data-Versioning: Overview of the fine grained mechanisms

This approach has several advantages:

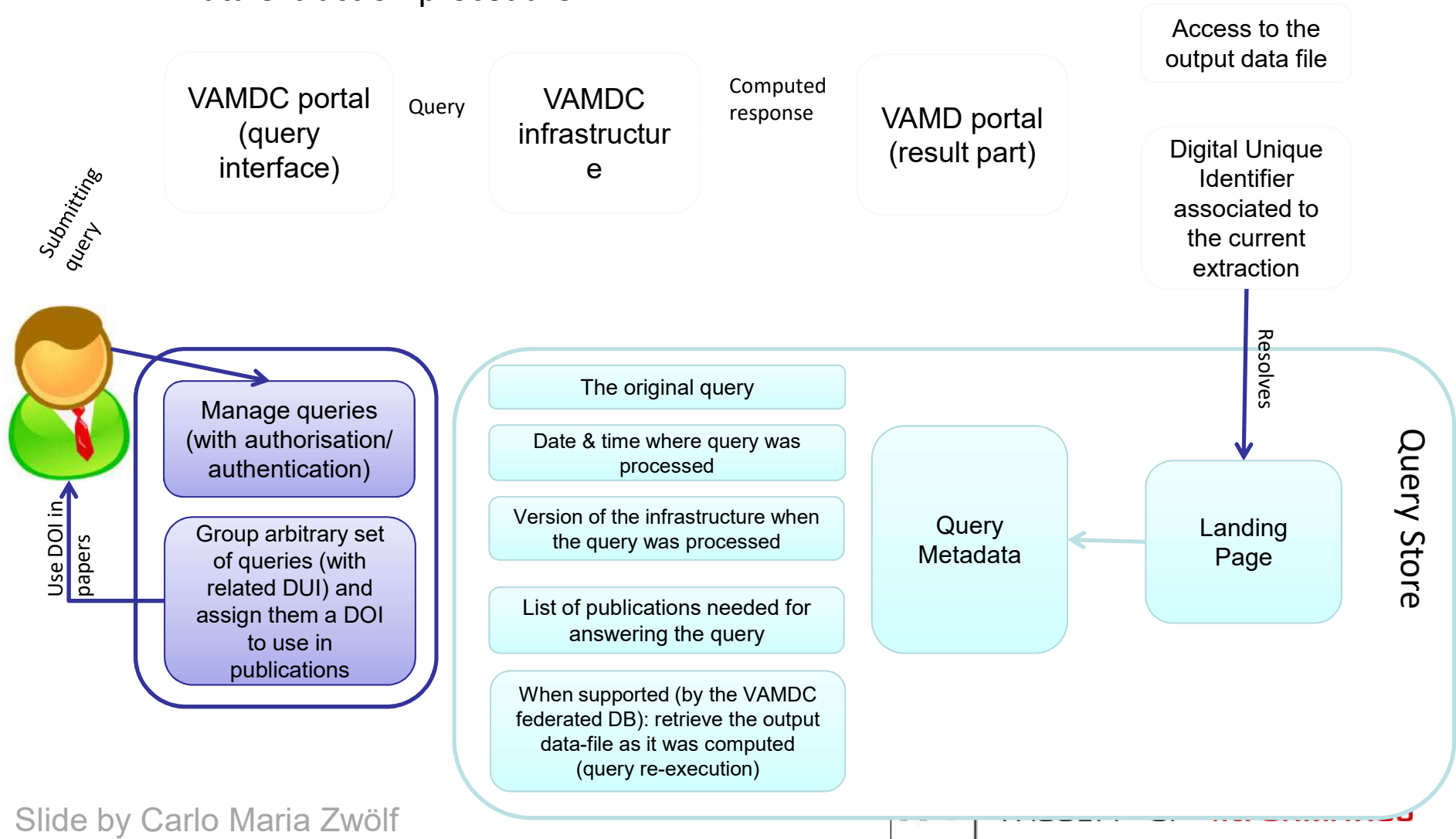
- It solves the data tagging granularity problem
- It is independent from what is considered a dataset
- The new files are compliant with old libraries & processing programs
 - We add a new feature, an overlay to the existing structure
 - We induce a structuration, without changing the structure

New model for datasets citation and extraction reproducibility in VAMDC,
C.M. Zwölf, N. Moreau, M.-L. Dubernet, *J. Mol. Spectrosc.* (2016),
<http://dx.doi.org/10.1016/j.jms.2016.04.009> Arxiv version:
<https://arxiv.org/abs/1606.00405>

Let us focus on the query store:

Sketching the functioning – From the final-user point of view:

Data extraction procedure



Let us focus on the query store:

The difficulty we have to cope with:

- Handle a query store in a distributed environment (RDA did not design it for these configurations)
- Integrate the query store with the existing VAMDC infrastructure

The implementation of the query store is the goal of a jointly collaboration between VAMDC and RDA-Europe

- Development started during spring 2016
- Final product released during 2017

Collaboration with Elsevier for embedding the VAMDC query store into the pages displaying the digital version of papers.

Designing technical solution for

- Paper / data linking at the paper submission (for authors)
- Paper / data linking at the paper display (for readers)



**Climate Change Centre Austria
(CCCA)**

Chris Schubert

chris.Schubert@ccca.ac.at

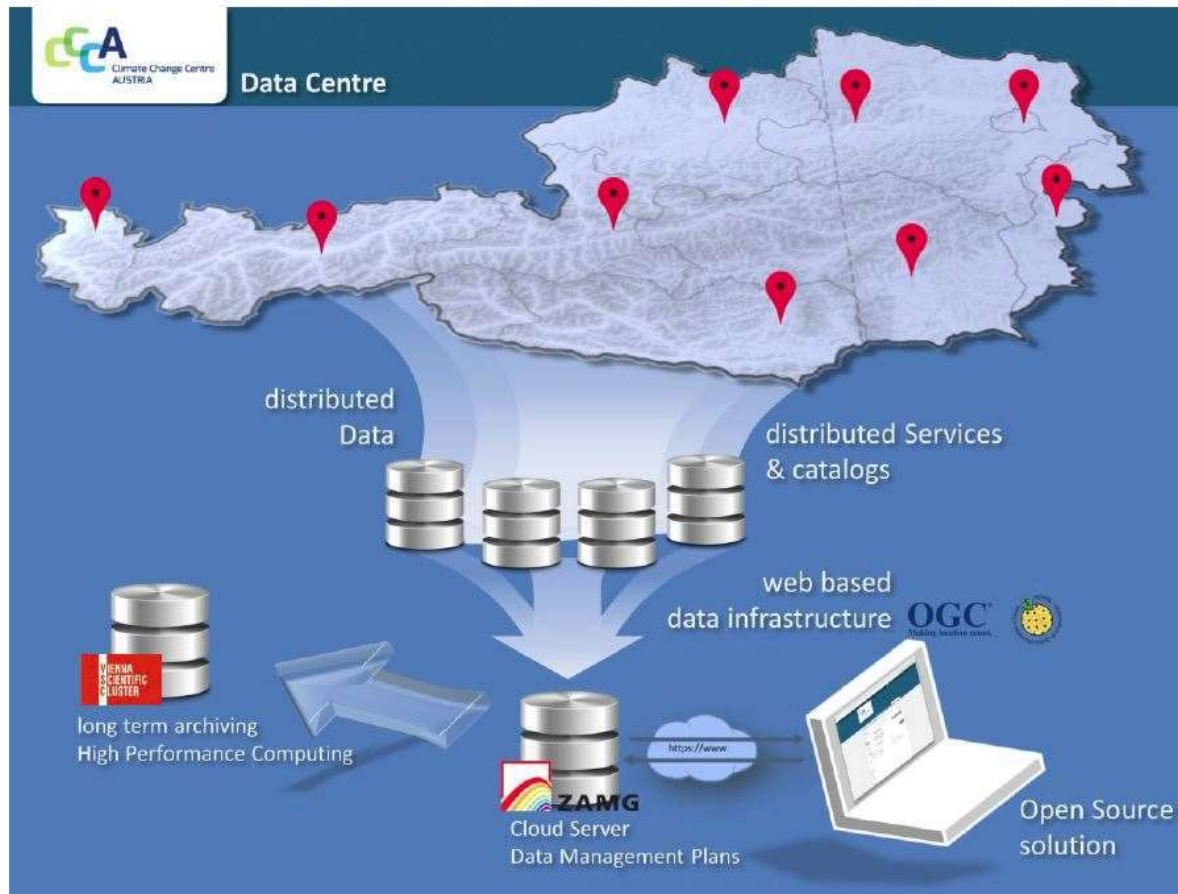
research data sharing without barriers

rd-alliance.org

Climate Change Centre Austria

- Climate research network for sustained, high-quality Austrian climate research.
- 28 members (11 universities, 13 non-university institutions, 4 supporting members)
- Structure: Coordination Office (Vienna, BOKU), Service Centre (Univ. Graz), Data Centre (ZAMG, Vienna)
- Service available at <http://data.ccca.ac.at>

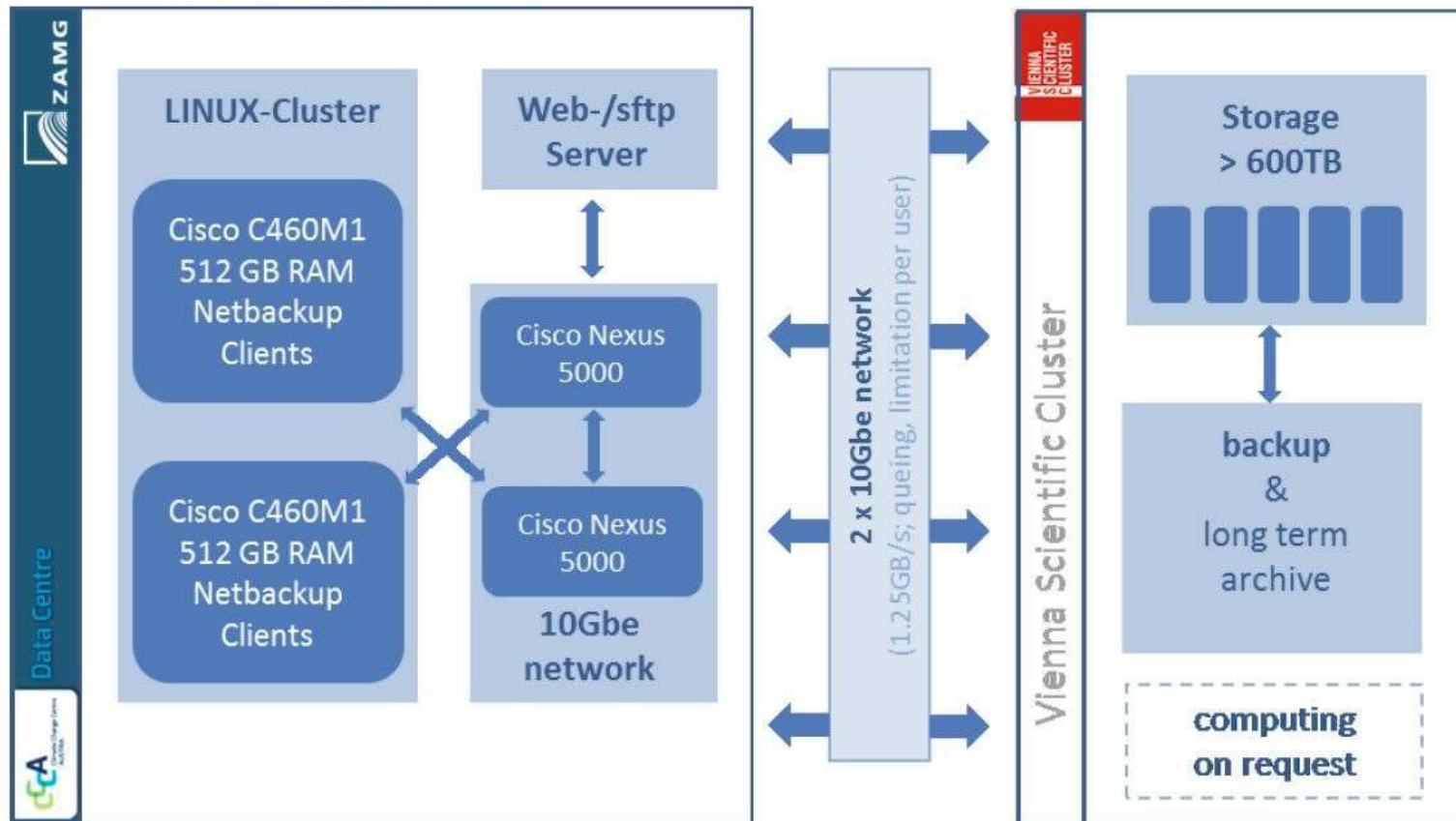




CCCA Data Centre

- › **provision** of climate-relevant information, data, algorithms, reports
- › **interoperable interfaces** to international portals, standards, legislation (e.g. INSPIRE)
- › conception for **long term archiving** of research data & repositories
- › capacity building, consultancy and **support for data sharing**

CCCA Data Centre Hardware





... a data portal among many others?

FEATURE No. 4 & 5

- handle® Service implemented to serve persistent identifier (PID) -> fundamental for DataCitation

hdl.handle.net/20.500.11756/7b9374de

Cite this resource:

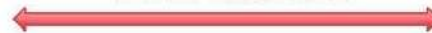
Using this data set or resource, you should cite this data set according to the given copyright conditions with following citation rules.

Hiebl et al. (2016). cdd-1961-2011-annual (Ver. 1). Retrieved from CCCA Data Centre: <https://hdl.handle.net/20.500.11756/fa338331>.
Access Date: February 22, 2017

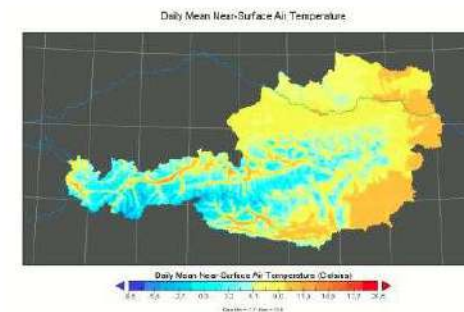
Your Publication



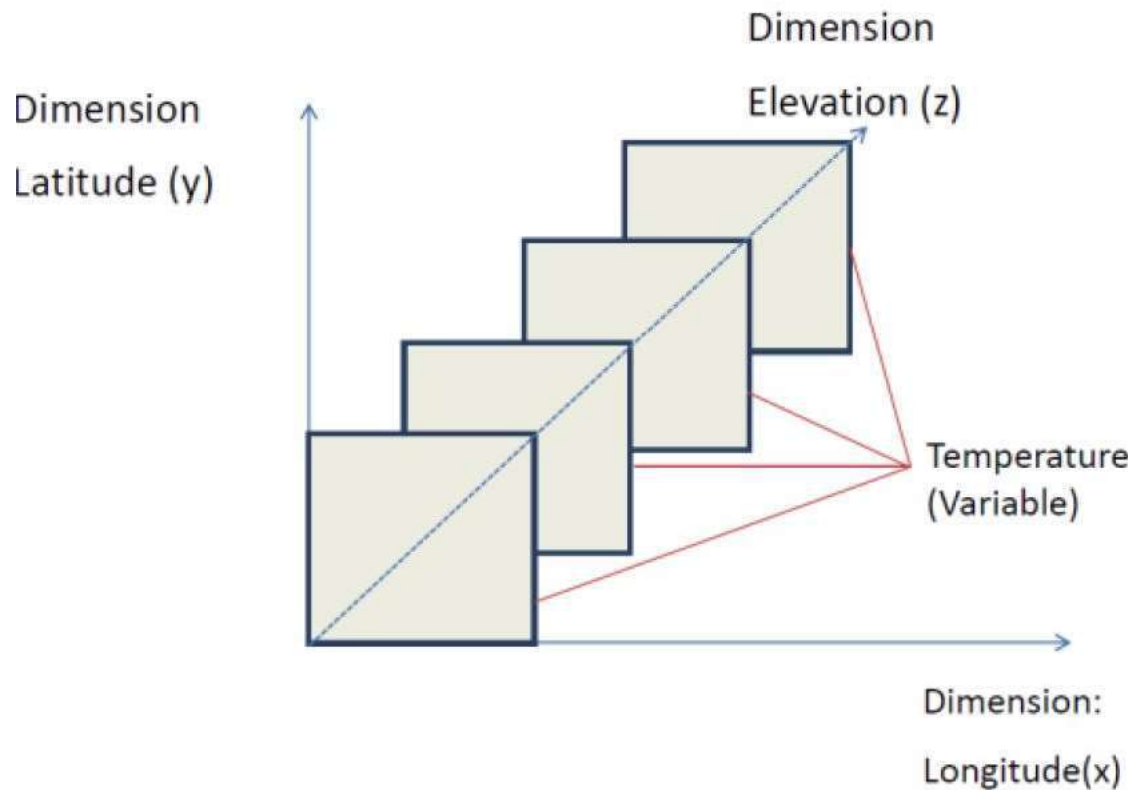
formal Data Citation



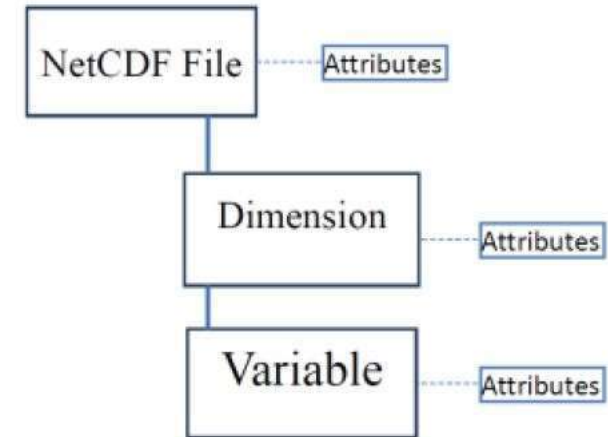
Your Data



- NetCDF Files:



modified and based on UCAR Unidata, www.unidata.ucar.edu/



```

* List of 3D
* $id : int 4
* $ndims : int 3
* $natts : int 7
* $unlimdimid : num 3
* $filename : chr "C:\VP8\Unidata\vol4\1779060-a951-11e1-b55e-e1515396a7,
* $varidindex: num [1:7] 0 0 0 1 2 3 4
* $writable : log FALSE
* $dim : List of 3
  ..$longitude:List of 8
  ...$name : chr "longitude"
  ...$len : int 720
  ...$units : log FALSE
  ...$id : int 1
  ...$dimaid : num 1
  ...$units : chr "degrees_east"
  ...$vals : num [1:720] 0.25 0.75 1.25 1.75 2.25 2.75 3.25 3.75 4.25 4.75 ...
  ...$create_dimvar: log TRUE
  ...$attr("class")= chr "dim.ncdf,"
  .....
* $name : num 4
* $var : List of 4
  ..$biomass_carbon_burning_varCF:List of 16
  ...$id : int 4
  ...$name : chr "biomass_carbon_burning_varCF"
  ...$ndims : int 3
  ...$natts : int 4
  ...$size : int [1:3] 720 270 4
  ...$prec : chr "float"
  ...$dimids : num [1:3] 1 2 3
  ...$units : chr "kg m-2"
  ...$longname : chr "biomass carbon burning"
  ...$dim : List()
  ...$dim : List of 3
  
```

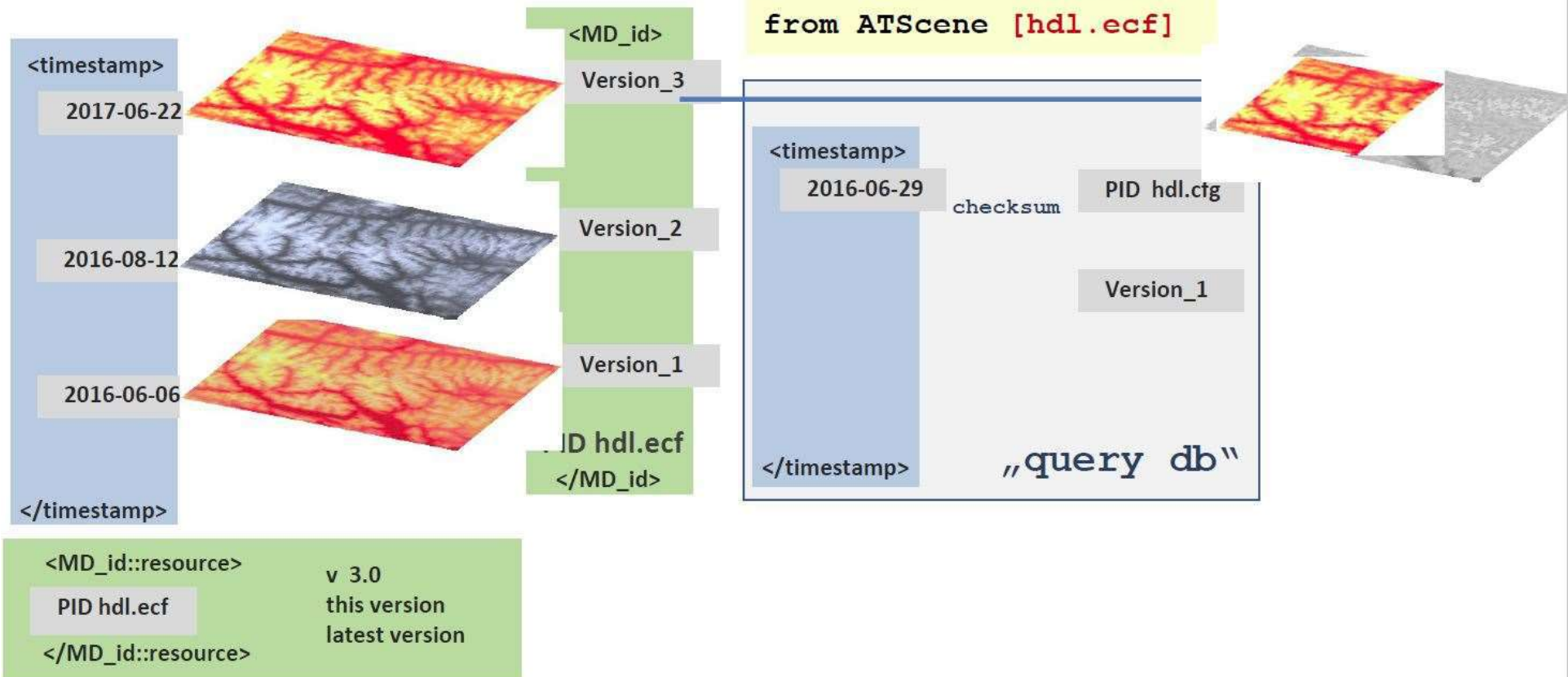
File Description

1st Dimension Description

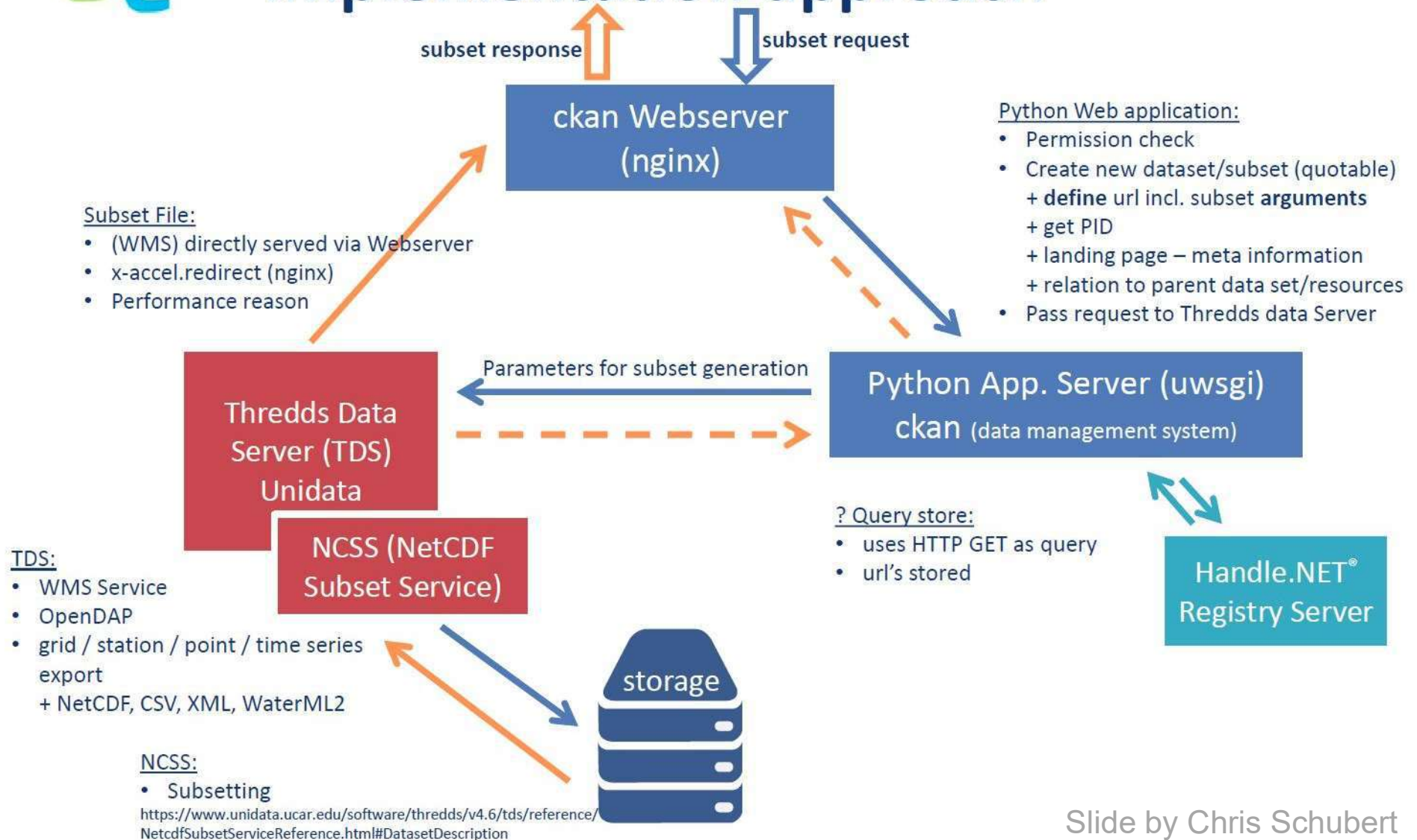
1st Variable Description

query var & lat/lon bbox

```
select var[tmax]
      lat/lon[48,14]
from ATScene [hdl.ecf]
```



Implementation approach



- for subsetting datasets
- uses **HTTP GET** as query in following scheme:
<http://{{host}}/{{context}}/{{service}}/{{dataset}}/{{dataset.html}} | {?query}>

Subsetting parameter used:

- **var** - names of our layer
- **north, south, east, west** - for the geographical extend, the bounding box
- **time_start, time_end, time_duration** - for time extend, limited only on 5 years interval
- **accept** - specify the returned format

All "http get" stored as url in our ckan data store

PID:

hdl.handle.net/20.500.11756/93887ecf

https://data.ccca.ac.at/tds_proxy/ncss/1dba52b2-4fd0-4fa1-a3ac-cfb0b94a7670?north=47.731688225506999&west=9.021605998277664&accept=netCDF&var=tas&east=12.031859904527664&south=46.77724203092812

Outline

-
- Why should we want to cite data?
 - What are the challenges in data identification and citation?
 - How should we do it, according to the RDA WG?
 - Who is doing it so far, and how?
 - **Summary**
-

Summary

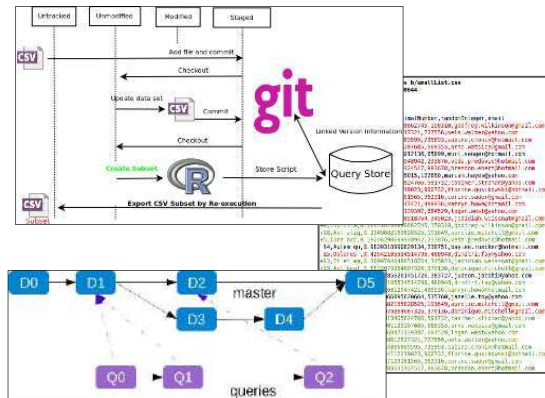
- Data citation essential for **solid** and **efficient** science
(but not just for science!)
- It is more than just giving credit
- Human-readable and machine-actionable
- RDA recommendations
 - Time-stamp and version data if it is evolving
 - Provide PIDs to arbitrary subsets via selection mechanism (“query”)
(rather than statically assigned PIDs to pre-defined subsets)
- 2 PIDs:
 - for evolving intellectual object
 - for precise, static subset

Benefits

- **Precisely identify any arbitrary subset of data**
- Principles applicable to all types of data
- Straightforward to implement in most settings
- Optimizations for high-volume / very dynamic data possible
- Transparent for the analyst / data scientist
- Reduces documentation effort for analysts / data scientist
- Reduces data management complexity for data centre
- Increases traceability of results, **trust**



Thank you!



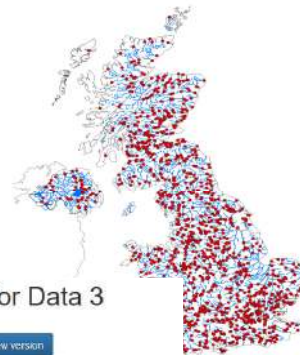
Label	Data	Real-time	Problem/Status	Method	Research Objective	Open
Repeat						Determinism
Param. Sweep	*					Robustness / Sensitivity
Generalize	+					Applicability across different
Port						Portability across platforms
Re-code	-					Correctness of implementations
Validate	+	+	+	+		Correctness of hypothesis, different approach
Re-use						Apply code in different settings, Repurpose
Independent & (orthogonal)						Suppleness of information, independent verification

Final Data Set	# WFs	
Disabled processors (WSDL services)	180	
Not executable in test environment	6	
Final Data Set	731	
Processor	# WFs	% WFs
Not terminated >48hours	6	0.8
Execution failed	384	52.5
Execution successful	341	46.6

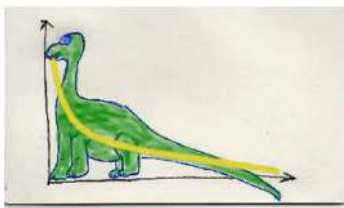
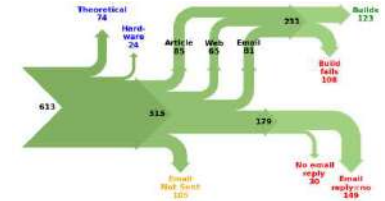
Diagram showing data flow: RDC table (c1, c2, c3) + sys_period triggers RDC_hist_table* (c1, c2, c3, sys_po). A note indicates '*stores history of data changes'.

Thanks!

<https://rd-alliance.org/working-groups/data-citation-wg.html>



DC¹
Data Citation Principles



Editing data for Data 3

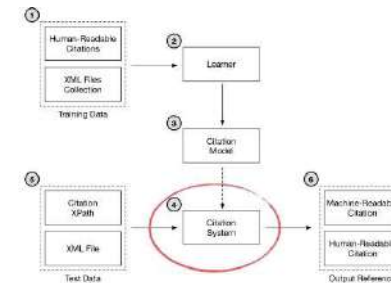
Back to Dataset Versions

Show changes Save to a new version

1 UPDATE 2000_1_test SET 'SiteID' = 'Stevensville Brook' WHERE db_table_pk=30
2 DELETE FROM 2000_1_test where db_table_pk=30
3 DELETE FROM 2000_1_test where db_table_pk=35

Actions	SiteID	LabID	Date	MeanDensity	Mean
	Stevensville Brook	2000.187	0000-00-00	46.44322354	39.0
	Winhall River	2011.081	2011-10-07	201	47.5
	Winhall River	2012.080	2012-09-27	1981	52.0
	Winhall	2013.170	2013-10-15	1002	30.0

2010



FAIR principles

Tomasz Miksa

TU Wien & SBA Research

tmiksa@sba-research.org

Agenda

- Introduction
- FAIR principles in detail
 - Important concepts underpinning the principles
 - Persistent identifiers, Metadata, Vocabularies, etc.
- FAIR assessment
- FAIR Digital Object
- Summary
- Literature and useful resources

INTRODUCTION

Hans Rosling and Data Science

Talk held in 2006

16 years later the problems are not solved

- (But a lot is going on to change this)

The TED logo is displayed in a bold, red, sans-serif font.

https://www.ted.com/talks/hans_rosling_the_best_stats_you_ve_ever_seen#t-1144167

Hans Rosling

“Because the **data is hidden down in the databases**. And the public is there, and the internet is there, but we have still not used it effectively.”

“There are some web pages like this,(...), but people **put prices** on them, **stupid passwords** and **boring statistics**.”

Hans Rosling

“Some countries accept that their databases can go out on the world. But what we really need is, of course, a search function, a search function where **we can copy the data up to a searchable format** and get it out in the world.”

“The **publicly funded data** is down here. (...) One of the crucial points is to **make them searchable**, and then people can use the different design tools to animate it there.”

Variety of solutions

In response to these needs many solutions were proposed and are being implemented

- **open access** to scientific publications and data
- research **data repositories** to host the data
- **persistent identifiers** to locate the data
- **data management plans**
- **FAIR principles**
- ...

Simplified view on FAIR

- Superficial
- For non-tech people

Discussed in the introductory lecture!

Findable – simplified examples

✓ Data repository
 ✗ Personal website

informatics

Accessible – simplified examples

✓ Restricted access, but a clear way to request access
 ✗

informatics

Interoperable – simplified examples

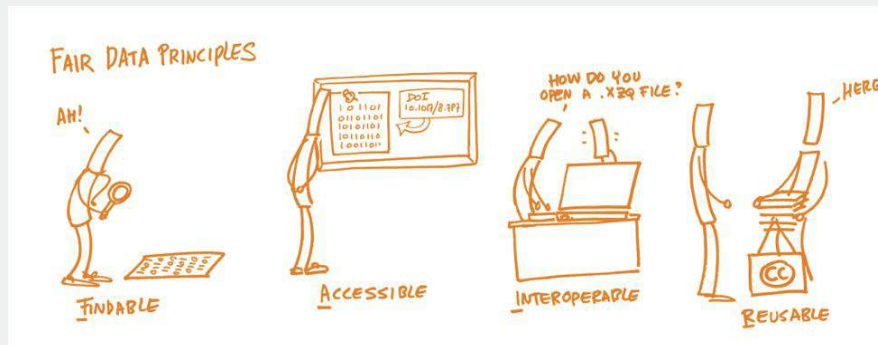
<p>Yes</p> <ul style="list-style-type: none"> • XML following known XSD Schema • MP3 for audio recordings • Data model using common vocabularies ✓	<p>No</p> <ul style="list-style-type: none"> • Custom XML without any documentation • M4P (Apple) for audio recordings • Custom fields in data model with poor documentation ✗
--	--

informatics

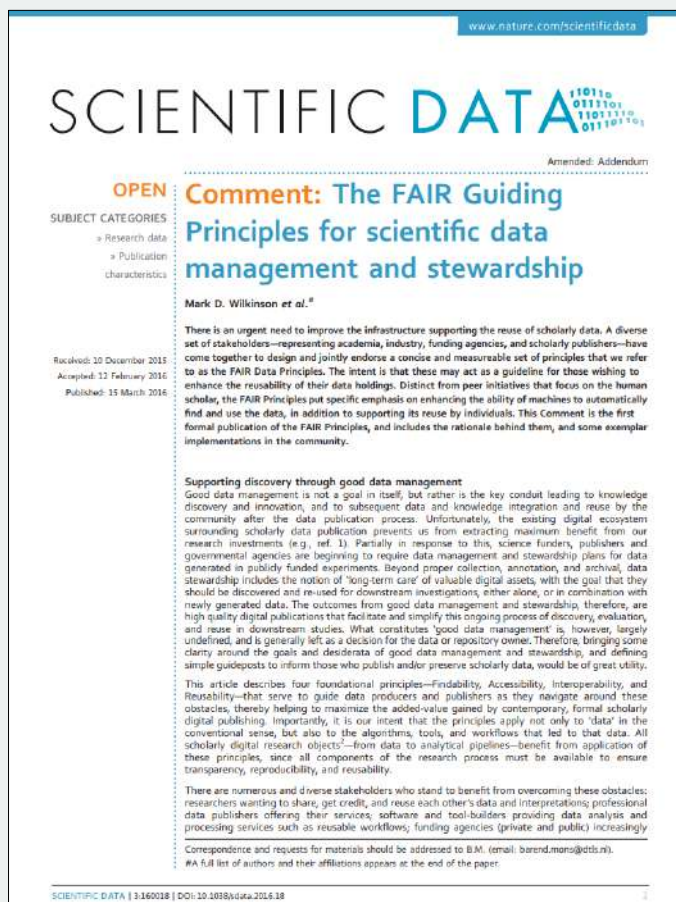
Reusable – simplified examples

✓ Trusted source, permission to reuse, well defined meaning of terms used
 ✗ Provenance and permissions not clear

informatics



FAIR principles



www.nature.com/scientificdata

SCIENTIFIC DATA

Amended: Addendum

Comment: The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson et al.^a

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplar implementations in the community.

Supporting discovery through good data management

Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our research investments (e.g., ref. 1). Partially in response to this, science funders, publishers and governmental agencies are beginning to require data management and stewardship plans for data generated in publicly funded experiments. Beyond proper collection, annotation, and archival, data stewardship includes the notion of ‘long-term care’ of valuable digital assets, with the goal that they should be discovered and re-used for downstream investigators, either alone, or in combination with newly generated data. The outcomes from good data management and stewardship, therefore, are high quality digital publications that facilitate and simplify this ongoing process of discovery, evaluation, and reuse in downstream studies. What constitutes ‘good data management’ is, however, largely undefined, and is generally left as a decision for the data or repository owner. Therefore, bringing some clarity around the goals and desiderata of good data management and stewardship, and defining simple guideposts to inform those who publish and/or preserve scholarly data, would be of great utility.

This article describes four foundational principles—Findability, Accessibility, Interoperability, and Reusability—that serve to guide data producers and publishers as they navigate around these obstacles, thereby helping to maximize the added-value gained by contemporary, formal scholarly digital publishing. Importantly, it is our intent that the principles apply not only to ‘data’ in the conventional sense, but also to the algorithms, tools, and workflows that led to that data. All scholarly digital research objects—from data to analytical pipelines—benefit from application of these principles, since all components of the research process must be available to ensure transparency, reproducibility, and reusability.

There are numerous and diverse stakeholders who stand to benefit from overcoming these obstacles: researchers wanting to share, get credit, and reuse each other’s data and interpretations; professional data publishers offering their services; software and tool-builders providing data analysis and processing services such as reusable workflows; funding agencies (private and public) increasingly.

Correspondence and requests for materials should be addressed to B.M. (email: barend.mon@edits.nl).
A full list of authors and their affiliations appears at the end of the paper.

SCIENTIFIC DATA | 3:160018 | DOI:10.1038/sdata.2016.18

<https://www.nature.com/articles/sdata201618>



GO FAIR

FAIR Principles Implementation Networks News Events Resources About GO FAIR Q

FAIR Principles

Home » FAIR Principles

- FAIR Principles
 - F1: (Meta) data are assigned globally unique and persistent identifiers
 - F2: Data are described with rich metadata
 - F3: Metadata clearly and explicitly include the identifier of the data they describe
 - F4: (Meta)data are registered or indexed in a searchable resource
 - A1: (Meta)data are retrievable by their identifier using a standardised communication protocol
 - A1.1: The protocol is open, free and universally implementable
 - A1.2: The protocol allows for an authentication and authorisation where necessary
 - A2: Metadata should be

In 2016, the **FAIR Guiding Principles for scientific data management and stewardship** were published in *Scientific Data*. The authors intended to provide guidelines to improve the Findability, Accessibility, Interoperability, and Reuse of digital assets. The principles emphasise machine-actionability (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention) because humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data.

A practical ‘how to’ guidance to go FAIR can be found in the **Three-point FAIRification Framework**.

Findable

The first step in (re)using data is to find them. Metadata and data should be easy to find for both humans and computers. Machine-readable metadata are essential for automatic discovery of datasets and services, so this is an essential component of the **FAIRification process**.

F1. (Meta)data are assigned a globally unique and persistent identifier

F2. Data are described with rich metadata (defined by R1 below)

F3. Metadata clearly and explicitly include the identifier of the data they describe

F4. (Meta)data are registered or indexed in a searchable resource

Accessible

Once the user finds the required data, she/he needs to know how can they be accessed, possibly including authentication and authorisation.

A1. (Meta)data are retrievable by their identifier using a standardised communications protocol

<https://www.go-fair.org/fair-principles/>

MACHINE-ACTIONABILITY

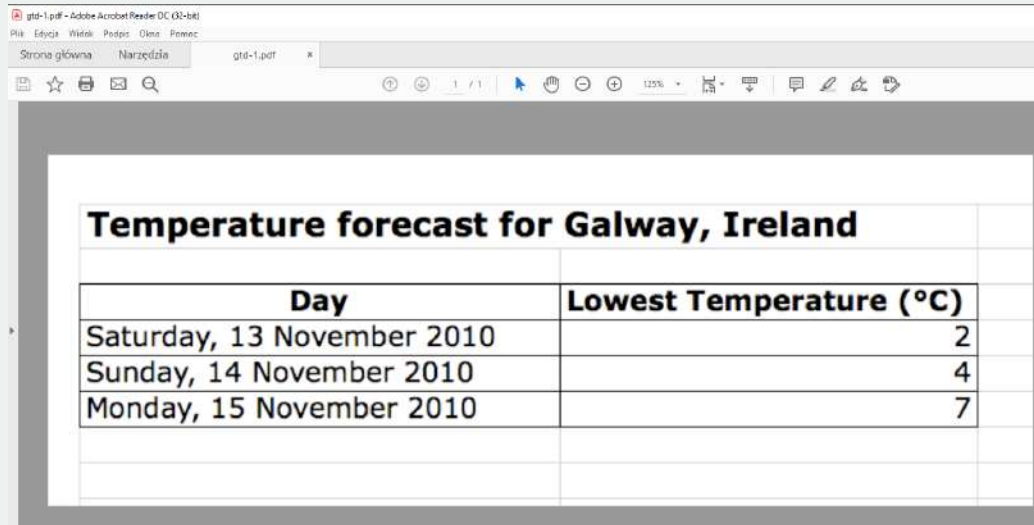
Machine actionability

“the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention” <https://www.go-fair.org/fair-principles>

“information that is structured in a consistent way so that machines, or computers, can be programmed against the structure.” <https://ddialliance.org/taxonomy/term/198>

Machine-actionability is core to each of the FAIR principles

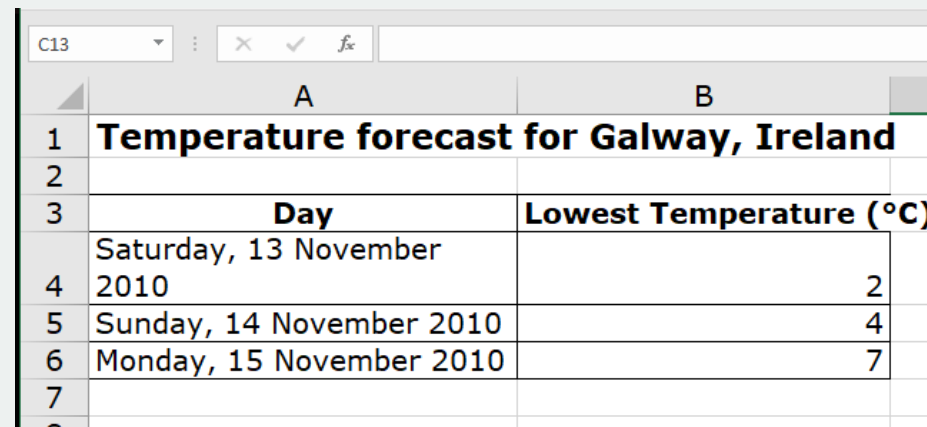
Machine-actionability - example



The screenshot shows a PDF document titled "gtd-1.pdf" in Adobe Acrobat Reader. The document contains a table with the following data:

Temperature forecast for Galway, Ireland	
Day	Lowest Temperature (°C)
Saturday, 13 November 2010	2
Sunday, 14 November 2010	4
Monday, 15 November 2010	7

Not machine-actionable



The screenshot shows a spreadsheet application with the same data as the PDF, but in a machine-actionable format. The data is organized as follows:

Temperature forecast for Galway, Ireland	
Day	Lowest Temperature (°C)
Saturday, 13 November 2010	2
Sunday, 14 November 2010	4
Monday, 15 November 2010	7

Machine-actionable

Machine-actionability – example (Linked Open Data)

```
Temperature forecast for Galway, Ireland
https://5stardata.info/en/examples/gtd-5/
  → dc:title → "Temperature forecast for Galway, Ireland"
  → html:stylesheet → https://5stardata.info/css/style.css
  → dcterms:title → "Temperature forecast for Galway, Ireland"
  → dcterms:created → "2012-01-22"^^xsd:date
  → dcterms:creator → http://mhausenblas.info/#i
  → dcterms:modified → "2015-08-31"^^xsd:date
  → dcterms:contributor → http://jayg.me/
  → dcterms:license → http://creativecommons.org/publicdomain/zero/1.0/

https://5stardata.info/en/examples/gtd-5/#Galway
  → rdf:type → meteo:Place
  → owl:sameAs → dbpedia:Galway
  → meteo:forecast → https://5stardata.info/en/examples/gtd-5/#forecast20101113, https://5stardata.info/en/examples/gtd-5/#forecast20101114, https://5stardata.info/en/examples/gtd-5/#forecast20101115

https://5stardata.info/en/examples/gtd-5/#forecast20101113
  → meteo:predicted → "2010-11-13T00:00:00Z"^^xsd:dateTime
  → meteo:temperature → https://5stardata.info/en/examples/gtd-5/#temp20101113
  ← is meteo:forecast of ← https://5stardata.info/en/examples/gtd-5/#Galway

https://5stardata.info/en/examples/gtd-5/#forecast20101114
  → meteo:predicted → "2010-11-14T00:00:00Z"^^xsd:dateTime
  → meteo:temperature → https://5stardata.info/en/examples/gtd-5/#temp20101114
  ← is meteo:forecast of ← https://5stardata.info/en/examples/gtd-5/#Galway

https://5stardata.info/en/examples/gtd-5/#forecast20101115
  → meteo:predicted → "2010-11-15T00:00:00Z"^^xsd:dateTime
  → meteo:temperature → https://5stardata.info/en/examples/gtd-5/#temp20101115
  ← is meteo:forecast of ← https://5stardata.info/en/examples/gtd-5/#Galway

https://5stardata.info/en/examples/gtd-5/#temp
  → rdfs:seeAlso → dbpedia:Temperature
  → owl:sameAs → dbpedia:Celsius

https://5stardata.info/en/examples/gtd-5/#temp20101113
  → meteo:celsius → "2"^^xsd:decimal
  ← is meteo:temperature of ← https://5stardata.info/en/examples/gtd-5/#forecast20101113

https://5stardata.info/en/examples/gtd-5/#temp20101114
  → meteo:celsius → "4"^^xsd:decimal
  ← is meteo:temperature of ← https://5stardata.info/en/examples/gtd-5/#forecast20101114

https://5stardata.info/en/examples/gtd-5/#temp20101115
  → meteo:celsius → "7"^^xsd:decimal
  ← is meteo:temperature of ← https://5stardata.info/en/examples/gtd-5/#forecast20101115
```

Unit definition and link to common definition

Values, types, link to forecast

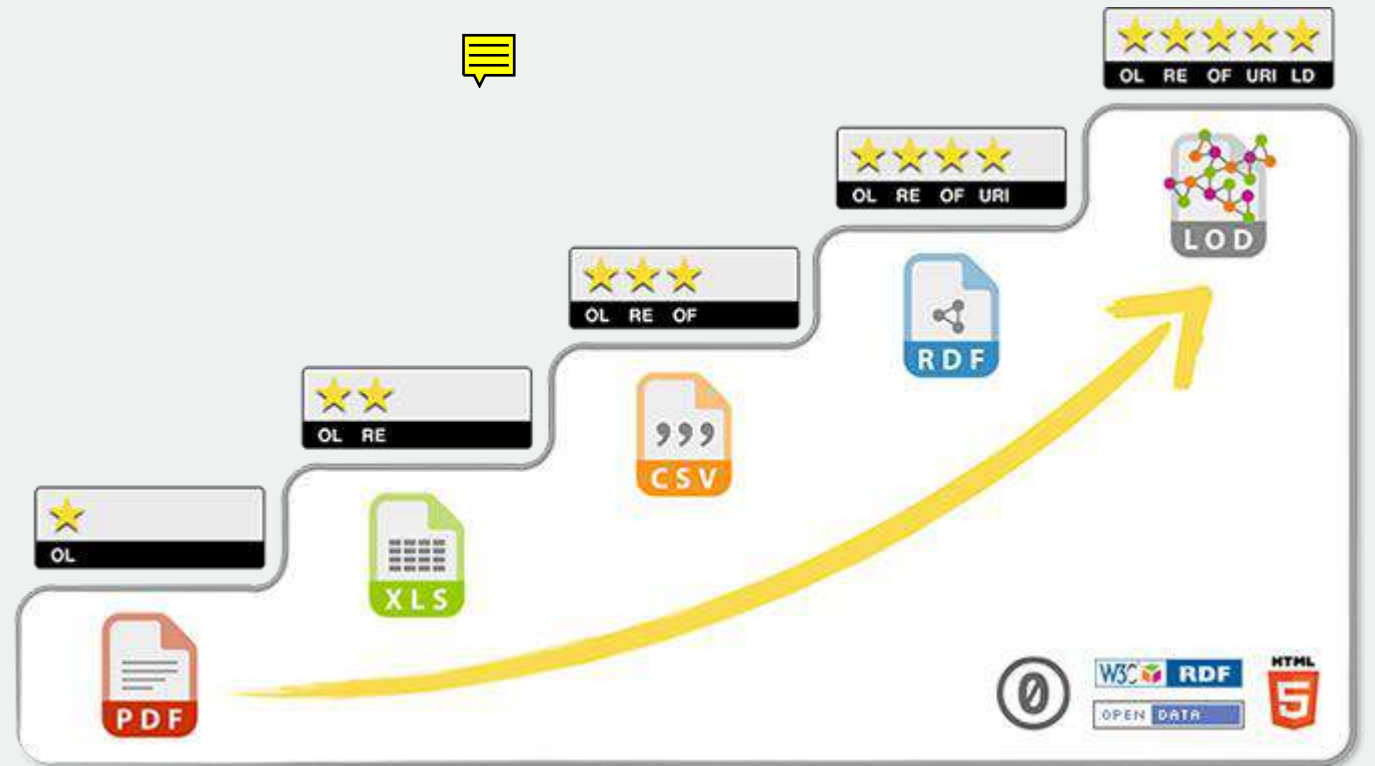
<http://graphite.ecs.soton.ac.uk/browser/?uri=http://5stardata.info/en/examples/gtd-5/>

Machine actionability – different shades

5-star model shows importance and benefits of

- machine-actionability
- open data
- semantic modelling

- ★ make your stuff available on the Web (whatever format) under an open license¹
- ★★ make it available as structured data (e.g., Excel instead of image scan of a table)²
- ★★★ make it available in a non-proprietary open format (e.g., CSV instead of Excel)³
- ★★★★ use URIs to denote things, so that people can point at your stuff⁴
- ★★★★★ link your data to other data to provide context⁵



Explore examples on the website to learn more on differences between each level!

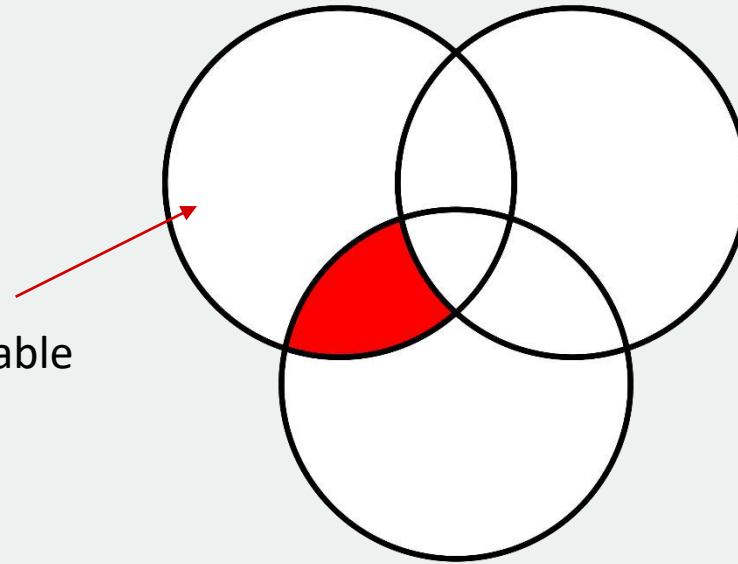
Machine actionability and FAIR

- Machine-actionability is core to each of the FAIR principles
- The more machine-actionable data is, the better it is
- FAIR does not require data to be open
 - 5-star model suggests openness – don't confuse those two!

Machine actionability and FAIR (how people often see it)

FAIR

Open



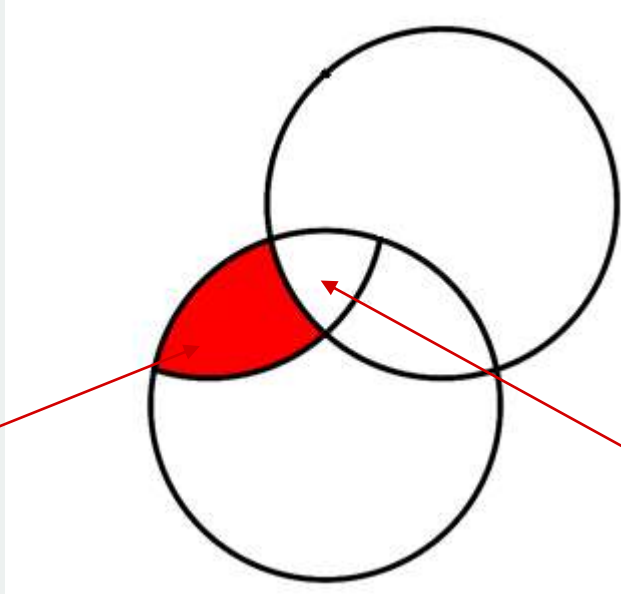
FAIR and not machine-actionable
(no such thing!)

Machine-actionable

Machine actionability and FAIR (how it really is)

FAIR

Open



FAIR must always be machine-actionable
and does not have to be open

FAIR must be machine-actionable
and can be open

Machine-actionable

FAIR IN DETAIL

Findable

F1. (Meta)data are assigned a globally unique and **persistent identifier**

F2. Data are described with rich **metadata**

F3. Metadata clearly and explicitly include the identifier of the data they describe

F4. (Meta)data are registered or indexed in a **searchable resource**

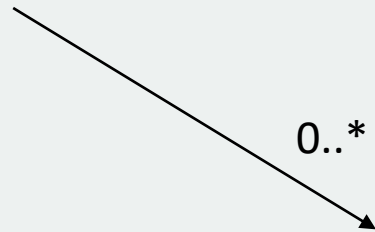
Persistent identifiers (F1)

Example

- A car has only one VIN (PID), but can have many number plates over its lifetime (URL)



VIN:	AZUSA1234567892222
● MODEL:	Awesome Car ●
DATE of MFG:	1970



Persistent Identifiers (F1)

Digital Object Identifier (DOI)

- Uniquely identify objects
- DOI assigned once
- Physical location of data can change



10.5281/zenodo.1068223

ORCID ID

- Unique person ID
- ORCID assigned once
- Person can change affiliations (jobs)



0000-0002-4929-7875

Persistent Identifiers - DOI

Unlike the URLs, DOIs are associated to objects and not to locations

URLs are unique, but not persistent

DOIs are never deleted

- if resource does not exist then a message is provided

Resolver service

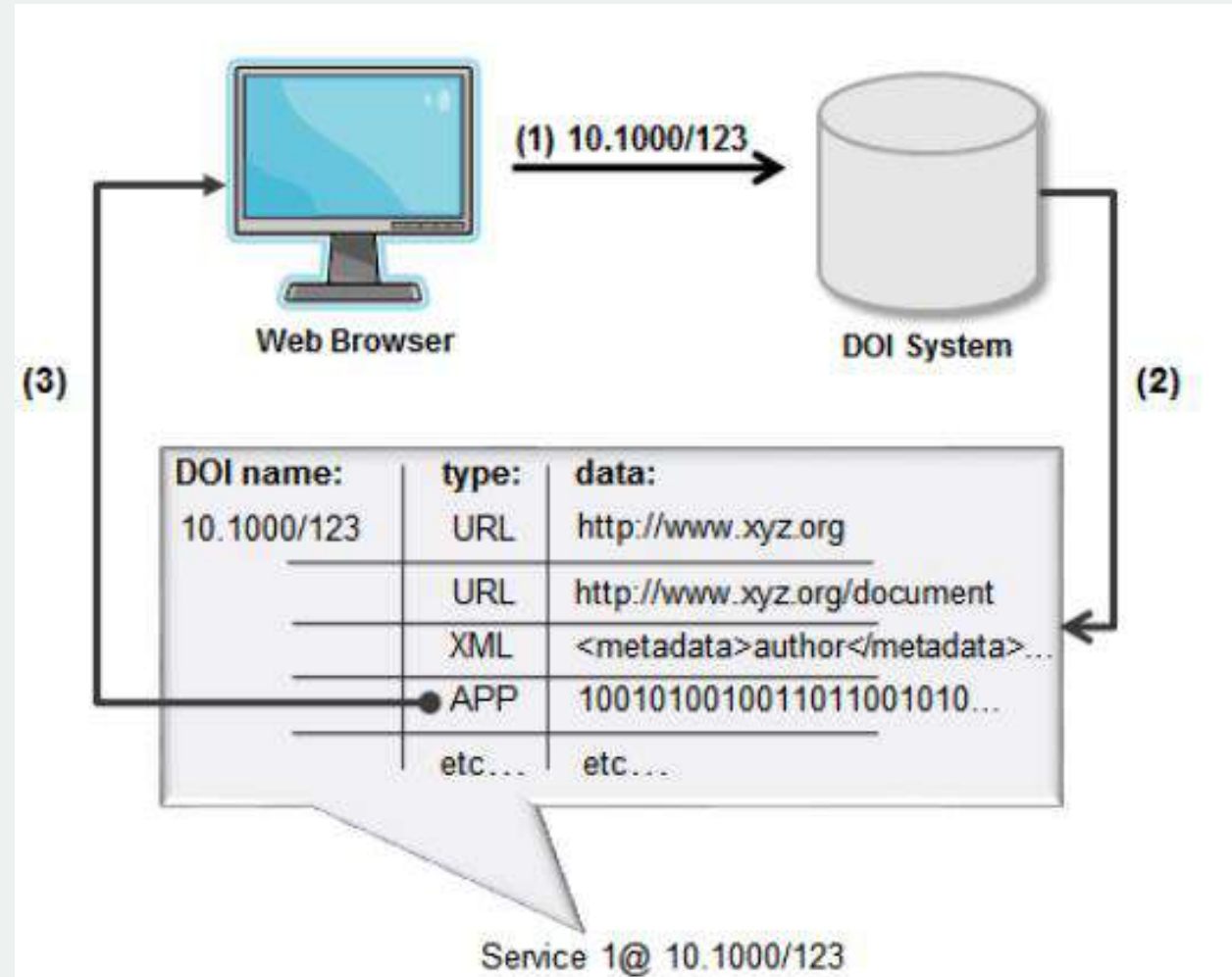
Metadata



`http://doi.org/ 10.4225 / 01/4F3DB08617645`

<code>http://doi.org/</code>	<code>10.4225</code>	<code>/ 01/4F3DB08617645</code>
resolver service	prefix (assigning body)	suffix (resource)

DOI – resolver service



https://www.doi.org/doi_handbook/3_Resolution.html

DOI example – assigned to publication

PLOS COMPUTATIONAL BIOLOGY

EDUCATION

Ten principles for machine-actionable data management plans

Tomasz Miksa¹*, Stephanie Simma², Daniel Mitchen³, Sarah Jones⁴

1 SBA Research & TU Wien, Vienna, Austria, **2** California Digital Library, University of California, Oakland, United States of America, **3** Data Science Institute, University of Virginia, Charlottesville, United States of America, **4** Digital Curation Centre, Glasgow, United Kingdom

* These authors contributed equally to this work.
* miksa@ifs.tuwien.ac.at

 Check for updates

OPEN ACCESS

Citation: Miksa T, Simma S, Mitchen D, Jones S (2019) Ten principles for machine-actionable data management plans. *PLoS Comput Biol* 15(3): e1006750. <https://doi.org/10.1371/journal.pcbi.1006750>

Published: March 28, 2019

Copyright: © 2019 Miksa et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: This research was carried out in the context of the Austrian COMET K1 program and publicly funded by the Austrian Research Promotion Agency (FFG) and the Vienna Business Agency (WAW). It was also supported by an NSF EXGER grant awarded to the California Digital Library (Award Number 1745575). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Abstract

Data management plans (DMPs) are documents accompanying research proposals and project outputs. DMPs are created as free-form text and describe the data and tools employed in scientific investigations. They are often seen as an administrative exercise and not as an integral part of research practice.

There is now widespread recognition that the DMP can have more thematic, machine-actionable richness with added value for all stakeholders: researchers, funders, repository managers, research administrators, data librarians, and others. The research community is moving toward a shared goal of making DMPs machine-actionable to improve the experience for all involved by exchanging information across research tools and systems and embedding DMPs in existing workflows. This will enable parts of the DMP to be automatically generated and shared, thus reducing administrative burdens and improving the quality of information within a DMP.

This paper presents 10 principles to put machine-actionable DMPs (maDMPs) into practice and realize their benefits. The principles contain specific actions that various stakeholders are already undertaking or should undertake in order to work together across research communities to achieve the larger aims of the principles themselves. We describe existing initiatives to highlight how much progress has already been made toward achieving the goals of maDMPs as well as a call to action for those who wish to get involved.

Introduction

Data management plans (DMPs) are documents accompanying research proposals. They describe the data that are used and produced during the course of research activities, where the data will be archived, which licenses and constraints apply, and to whom credit should be given. DMPs are awareness tools to help researchers manage their data and ensure that it will be of high quality, accessible, and reusable after the project has ended. DMPs are typically created manually, mostly by researchers using checklists and online questionnaires. They are required by funding bodies and institutions all over the world, e.g., the National Science

PLOS Computational Biology | <https://doi.org/10.1371/journal.pcbi.1006750> March 28, 2019 1 / 15

DOI example – assigned to code

The screenshot shows the GitHub repository page for 'helmuthb / dmp-exercise1'. The repository is at version 1.0.1, has 1 branch, and 3 tags. A commit by 'helmuthb' titled 'Corrected DOI link' is highlighted, with a green checkmark and the DOI '683c723' on 'Apr 22, 2019' with 4 commits. Below the commit list is a table of files:

File	Description	Time
data	First version with full data	2 years ago
src	First version with full data	2 years ago
.gitignore	First version with full data	2 years ago
Dockerfile	First version with full data	2 years ago
LICENSE	Initial commit	2 years ago
README.md	Corrected DOI link	2 years ago
Report.pdf	First version with full data	2 years ago

Below the file list, the 'README.md' content is shown, with a red box highlighting the DOI: `DOI 10.5281/zenodo.2648326`

The screenshot shows the Zenodo record page for 'US Wheat and Salzburg Middle-Aged Marriages - Data Experiment'. The record has 14 views and 5 downloads. It is associated with GitHub and OpenAIRE. The DOI '10.5281/zenodo.2648326' is highlighted with a red box. The record also shows a list of files and a table of versions.

File	Size
LICENSE	1.1 KB
README.md	232 Bytes
Report.pdf	246.6 KB
data	901 Bytes
metadata.yml	3.9 KB
processed	30.4 KB
us-wheat-sbg-marriages.csv	9.9 KB
wheat	1.3 MB
wheat_data-40_years.xls	20.1 KB
marriage-sbg-age-rsv	9.4 KB
src	9.4 KB
Report.html	4.1 KB
docker scripts	3.6 KB
install-stata-packages	126 Bytes
install-r-packages	30.4 KB
latex-pa-kages	4.1 KB
fenchel.tex.gz	3.6 KB
langage.tex.gz	1.3 KB
packages	1.3 KB
Diagrammi_1_9_0.tex.gz	626.0 KB
DiagrammiLang_0.1.tar.gz	626.0 KB

Version	Created
version 1.0.1	Apr 22, 2019
10.5281/zenodo.2648326	
Version 1.0	Apr 22, 2019
10.5281/zenodo.2648326	
Version 0.1	Apr 22, 2019
10.5281/zenodo.2648326	

DOI example - assigned to data

November 27, 2020 | Version 1.0 Dataset Embargoed

The Sentinel-1 Global Backscatter Model (S1GBM) - Mapping Earth's Land Surface with C-Band Microwaves

Bauer-Marschallinger, Bernhard ¹; Cao, Senmao ^{1,2}; Navacchi, Claudio ¹; Freeman, Vahid ^{1,3}; Reuß, Felix ¹; Geudtner, Dirk ⁴; Rommen, Björn ⁴; Vega, Francisco Ceba ⁴; Snoeij, Paul ⁵; Attema, Evert ⁴; Reimer, Christoph ²; Wagner, Wolfgang ^{1,2} [show affiliations](#)

Description
This dataset was generated by the Remote Sensing Group of the [TU Wien Department of Geodesy and Geoinformation \(https://mrs.geo.tuwien.ac.at/\)](https://mrs.geo.tuwien.ac.at/), within a dedicated project by the European Space Agency (ESA). Rights are reserved with ESA. Open use is granted under the CC BY-SA 4.0 license.

With this dataset publication, we open up a new perspective on Earth's land surface, providing a normalised microwave backscatter map from spaceborne Synthetic Aperture Radar (SAR) observations. The Sentinel-1 Global Backscatter Model (S1GBM) describes Earth for the period 2016-17 by the mean C-band radar cross section in VV- and VH-polarization at a 10 m sampling, giving a high-quality impression on surface- structures and -patterns.

At TU Wien, we processed 0.5 million Sentinel-1 scenes totaling 1.1 PB and performed semi-automatic quality curation and backscatter harmonisation related to orbit geometry effects. The overall mosaic quality excels (the few existing datasets, with minimised imprinting from orbit discontinuities and successful angle normalisation in large parts of the world. Supporting the designand verification of upcoming radar sensors, the obtained S1GBM data potentially also serve land cover classification and determination of vegetation and soil states, as well as water body mapping.

We invite developers from the broader user community to exploit this novel data resource and to integrate S1GBM parameters in models for various variables of land cover, soil composition, or vegetation structure.

Versions

Version 1.0
DOI: 10.48436/n2d1v-gqb91

Cite As

Bauer-Marschallinger, Bernhard et al. (2020). The Sentinel-1 Global Backscatter Model (S1GBM) - Mapping Earth's Land Surface with C-Band Microwaves (Version 1.0) [Dataset]. TU Data.
<https://doi.org/10.48436/n2d1v-gqb91>

Persistent Identifiers

More on persistent identifiers (PIDs) in the other lectures

Note: DOI and ORCID are not the only PIDs in use!

Findable (F2-F4)

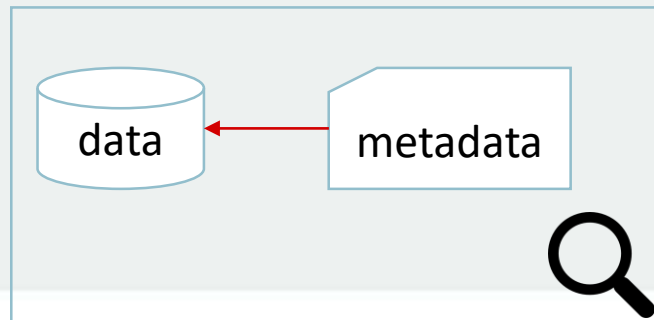
F2. Data are described with rich metadata



F3. Metadata clearly and explicitly include the identifier of the data they describe



F4. (Meta)data are registered or indexed in a searchable resource



Findable

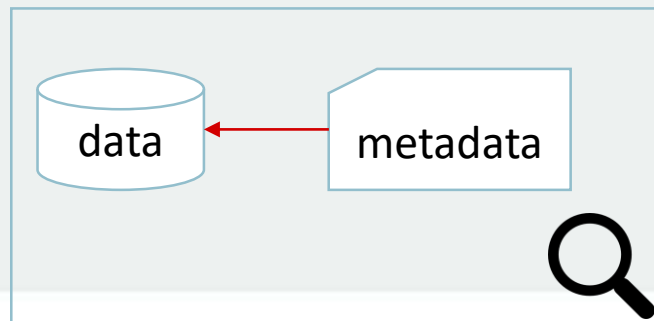
F2. Data are described with rich metadata



F3. Metadata clearly and explicitly include the identifier of the data they describe



F4. (Meta)data are registered or indexed in a searchable resource

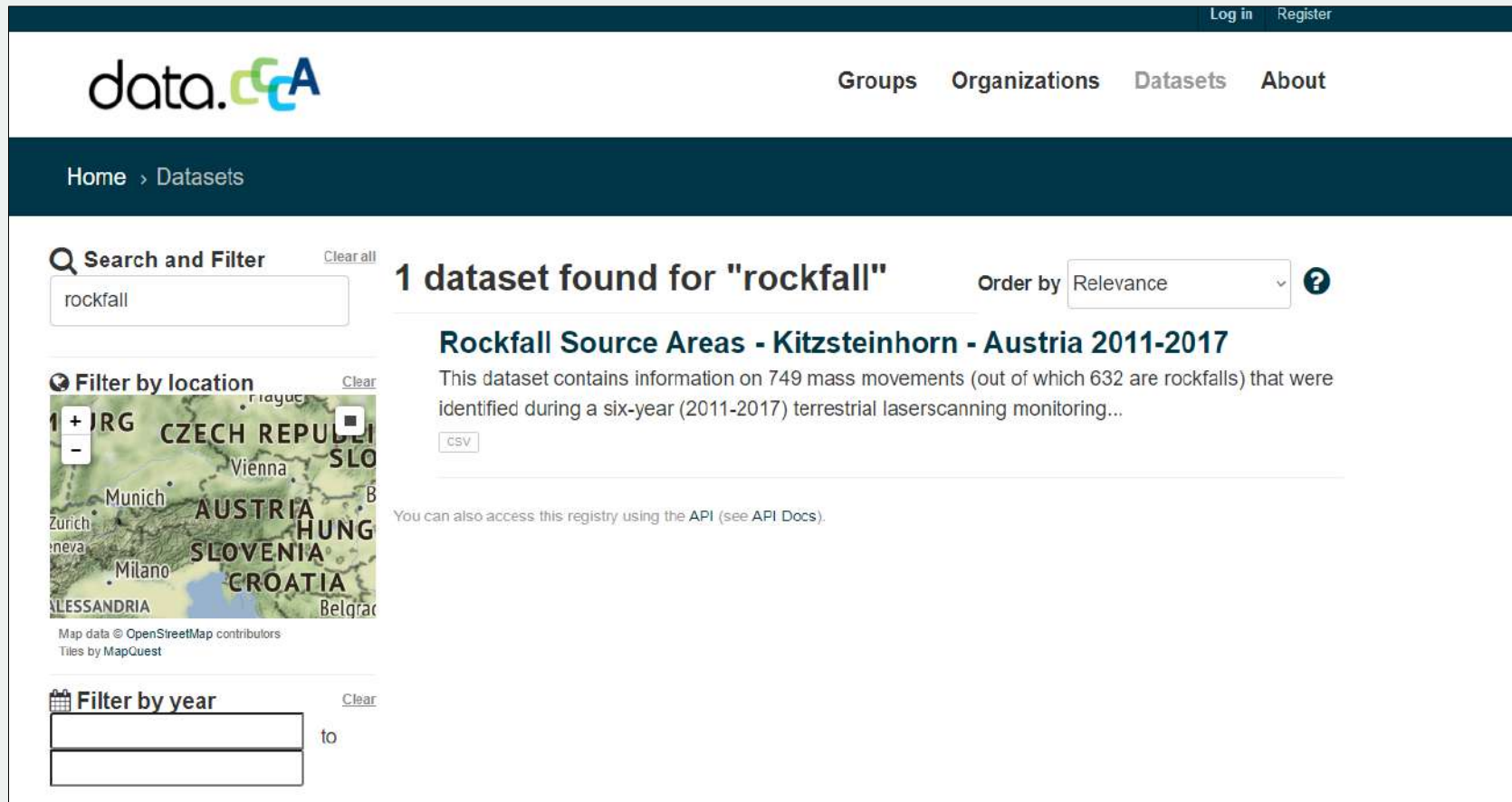


Would a TXT be ok instead?



/home/... ?
http:// ?
Handle ?

F4. (Meta)data are registered or indexed in a searchable resource



The screenshot shows the data.cca website interface. At the top right, there are links for "Log in" and "Register". The main navigation bar includes "Groups", "Organizations", "Datasets", and "About". Below this, a breadcrumb trail shows "Home > Datasets".

The search section is titled "Search and Filter" and contains a search box with the text "rockfall" and a "Clear all" link. To the right of the search box, it says "1 dataset found for 'rockfall'" and "Order by Relevance" with a dropdown menu and a help icon.

The search results display a dataset titled "Rockfall Source Areas - Kitzsteinhorn - Austria 2011-2017". The description states: "This dataset contains information on 749 mass movements (out of which 632 are rockfalls) that were identified during a six-year (2011-2017) terrestrial laserscanning monitoring...". There is a "CSV" download button.

Below the search box, there is a "Filter by location" section with a map of Central Europe. The map shows Austria, Czech Republic, Slovakia, Slovenia, Hungary, and Croatia. Major cities like Vienna, Munich, Zurich, and Milano are labeled. There are zoom controls and a "Clear" link.

At the bottom, there is a "Filter by year" section with two input boxes and a "to" label, and a "Clear" link.

More on the F4 in the lecture on repositories.

<https://data.cca.ac.at>

F3. Metadata clearly and explicitly include the identifier of the data they describe

The screenshot shows a dataset page on the 'data.cca' platform. The dataset title is 'Rockfall Source Areas - Kitzsteinhorn - Austria 2011-2017'. The description states it contains information on 749 mass movements (632 rockfalls) from 2011-2017. The page includes sections for 'Dataset Versions', 'Citation', 'Resources', and 'Dataset Metadata'. The 'Dataset Metadata' section has tabs for 'Contact', 'Basics', 'Keywords', 'Spatial', 'Time', 'Specifics', 'Quality', and 'Conformity'. The 'Basic Information about this dataset' section shows the 'Dataset Locator' as 'https://hdl.handle.net/20.500.11756/70ef62e8'. Annotations in red boxes and text highlight the 'Dataset' resource, the 'Export Metadata' button, and the 'Dataset Locator' URL.

Log in Register

data.cca Groups Organizations Datasets About

Home > Organizations > GEORESEARCH ... > Rockfall Source Areas - ...

DATASET

Rockfall Source Areas - Kitzsteinhorn - Austria 2011-2017

Followers: 0
Views: 22

Published by: GEORESEARCH Forschungsgesellschaft mbH License: Creative Commons Attribution - Share-Alike (CC-BY-SA)

This dataset contains information on 749 mass movements (out of which 632 are rockfalls) that were identified during a six-year (2011-2017) terrestrial laserscanning monitoring at the Kitzsteinhorn, Hohe Tauern Range, Austria. The data documents the significant impact that retreating glaciers have on rockfall occurrence in two deglaciating cirques. The dataset includes: mass movement volume, substrate type, failure depth, height of source area above the glacier surface, slope angle/aspect of source area. An extensive analysis and interpretation of the dataset can be found in two research papers published in the open-access journal "Earth Surface Dynamics" (Hartmeyer et al. 2020). Funding information: Data acquisition was co-funded by the Austrian Academy of Sciences (ÖAW) (Project 'GlacierRocks') and the Austrian Research Promotion Agency (FFG) (Project 'MOREXPRT').

Dataset Versions Citation

Resources

Dataset

Rockfall Source Areas, Kitzsteinhorn, Austria ...
This dataset contains information on 749 mass movements (out of which 632 are ...)

Explore

Dataset Metadata

Contact Basics Keywords Spatial Time Specifics Quality Conformity

Machine-readable metadata

Export Metadata

Basic Information about this dataset


Identifier

Dataset Locator - <https://hdl.handle.net/20.500.11756/70ef62e8>

<https://hdl.handle.net/20.500.11756/70ef62e8>

F2. Data are described with rich metadata

Resources

 **Rockfall Source Areas, Kitzsteinhorn, Austria ...** [Explore](#)

This dataset contains information on 749 mass movements (out of which 632 are...

Dataset Metadata

[Export Metadata](#)


Contact Basics **Keywords** Spatial Time Specifics Quality Conformity

Basic Information about this dataset

Dataset	https://hdl.handle.net/20.500.11756/70ef62e8
Locator - URI	
Abstract	This dataset contains information on 749 mass movements (out of which 632 are rockfalls) that were identified during a six-year (2011-2017) terrestrial laserscanning monitoring at the Kitzsteinhorn, Hohe Tauern Range, Austria. The data documents the significant impact that retreating glaciers have on rockfall occurrence in two deglaciating cirques. The dataset includes: mass movement volume, substrate type, failure depth, height of source area above the glacier surface, slope angle/aspect of source area. An extensive analysis and interpretation of the dataset can be found in two research papers published in the open-access journal "Earth Surface Dynamics" (Hartmeyer et al. 2020). Funding information: Data acquisition was co-funded by the Austrian Academy of Sciences (ÖAW) (Project 'GlacierRocks') and the Austrian Research Promotion Agency (FFG) (Project 'MOREXPART').
Metadata Language	English
License	cc-by-sa
Visibility	public
Use Limitation	no limitation

F2. Data are described with rich metadata

Resources

 **Rockfall Source Areas, Kitzsteinhorn, Austria ...**
This dataset contains information on 749 mass movements (out of which 632 are...)

[Explore](#)


Dataset Metadata [Export Metadata](#)

Contact Basics **Keywords** **Spatial** Time Specifics Quality Conformity

Geographic Aspects of the Resources

Polygon

Dataset extent



Map data © OpenStreetMap contributors
Tiles by MapQuest

Coverage Kitzsteinhorn, Hohe Tauern Range, Austria

Metadata - Chemistry


NIH National Library of Medicine
National Center for Biotechnology Information

PubChem About Blog Submit Contact

COMPOUND SUMMARY

Water ethanol


PubChem CID 19096565



Structure 
2D 3D
[Find Similar Structures](#)

Molecular Formula C2H6O2

Synonyms ethanol water
ethanol-water
water ethanol
water-ethanol
EtOH water
[More...](#)

Molecular Weight 64.08 g/mol

Parent Compound  CID 702 (Ethanol)

Component Compounds  CID 702 (Ethanol)
 CID 962 (Water)

Dates Modify 2021-02-27 Create 2007-12-04

2 Names and Identifiers

2.1 Computed Descriptors

2.1.1 IUPAC Name

ethanol;hydrate

Computed by LexiChem 2.6.6 (PubChem release 2019.06.18)

[PubChem](#)

2.1.2 InChI

InChI = 1S/C2H6O.H2O/c1-2-3;/h3H,2H2,1H3;1H2

Computed by InChI 1.0.5 (PubChem release 2019.06.18)

[PubChem](#)

2.1.3 InChI Key

IDGUHHHQCSQLU-UHFFFAOYSA-N

Computed by InChI 1.0.5 (PubChem release 2019.06.18)

[PubChem](#)

2.1.4 Canonical SMILES

CCO.O

Computed by OEChem 2.1.5 (PubChem release 2019.06.18)

[PubChem](#)

2.2 Molecular Formula

C2H8O2

Computed by PubChem 2.1 (PubChem release 2019.06.18)

[PubChem](#)

3 Chemical and Physical Properties

3.1 Computed Properties

Property Name	Property Value
Molecular Weight	64.08 g/mol
Hydrogen Bond Donor Count	2
Hydrogen Bond Acceptor Count	2
Rotatable Bond Count	0
Exact Mass	64.052429 g/mol
Monoisotopic Mass	64.052429 g/mol
Topological Polar Surface Area	21.2 Å ²
Heavy Atom Count	4
Formal Charge	0
Complexity	2.8
Isotope Atom Count	0
Defined Atom Stereocenter Count	0
Undefined Atom Stereocenter Count	0
Defined Bond Stereocenter Count	0
Undefined Bond Stereocenter Count	0
Covalently-Bonded Unit Count	2
Compound Is Canonicalized	Yes

[PubChem](#)

Accessible

A1. (Meta)data are retrievable by their identifier using a standardised communications protocol

- A1.1 The protocol is open, free, and universally implementable
- A1.2 The protocol allows for an authentication and authorisation procedure, where necessary

A2. Metadata are accessible, even when the data are no longer available

A1.1 The protocol is open, free, and universally implementable

“Anyone with a computer and an internet connection can access at least the metadata”

HTTP

- Open – specification of the protocol is known to everyone
- Free – no need to pay to “use Internet”

Proprietary protocols

- evade

OSI model		
Layer	Name	Example protocols
7	Application Layer	HTTP, FTP, DNS, SNMP, Telnet
6	Presentation Layer	SSL, TLS
5	Session Layer	NetBIOS, PPTP
4	Transport Layer	TCP, UDP
3	Network Layer	IP, ARP, ICMP, IPSec
2	Data Link Layer	PPP, ATM, Ethernet
1	Physical Layer	Ethernet, USB, Bluetooth, IEEE802.11

A1.2 protocol allows for authentication and authorisation

Protected and private data can be FAIR

Possible types of access

- Open – everyone has access
- Shared or restricted – only a selected/ invited group of people can access
- Closed – only the owner has access



Accessible - example

The screenshot shows the Zenodo search results page. The top navigation bar is blue with the Zenodo logo, a search bar, and links for 'Upload' and 'Communities'. A user profile 'miksa@ifs.tuwien.ac.at' is visible in the top right. The main content area shows search results for 'Found 1738119 results.' with a pagination bar and sorting options. Two search results are displayed, each with a 'View' button. The first result is 'Desk-Research Analysis and Identification of SA and Training Tools' by Mateusz Macias, dated March 5, 2021, with 'Project deliverable' and 'Open Access' tags. The second result is 'An embedded device for indoor localization in BLE networks based on a reconfigurable antenna' by Luszczak, Przemyslaw, dated January 15, 2021, with 'Thesis' and 'Open Access' tags. On the left side, there are two filter panels: 'Access Right' and 'File Type'. The 'Access Right' panel is highlighted with a red box and contains the following options: 'Open (1699699)', 'Closed (32706)', 'Restricted (4520)', and 'Embargoed (1194)'. The 'File Type' panel contains options: 'Pdf (892059)', 'Jpg (361789)', 'Png (221819)', 'Html (85105)', and 'Zip (79205)'. The 'Restricted' option in the 'Access Right' panel is highlighted with a red box.

<https://zenodo.org/search?page=1&size=20&q=>

Accessible - example

The screenshot shows the Zenodo interface for a dataset. At the top, the Zenodo logo is on the left, and a search bar, 'Upload', and 'Communities' links are on the right. Below the header, the date 'March 5, 2021' is on the left, and two buttons, 'Dataset' and 'Restricted Access', are on the right. The main title is 'Phase-contrast X-ray tomography of free-breathing murine lungs'. Below the title, the authors are listed: Kian Shaker, Ilian Häggmark, Jakob Reichmann, Marie Arsenian-Henriksson, and Hans M. Hertz. The description follows: 'Full resolution phase-contrast X-ray tomographic dataset acquired of the lungs of a free-breathing mouse (NMRI nude mice, BomTac:NMRI-Foxn1tm, Taconic Biosciences, DK), weighing roughly 28 g. Datasize: 3851x3951 pixels per slice, stack of 1700 slices, 16-bit, .tif. Voxel-size: 8.25x8.25x8.25 micrometer. A sample slice is available for download, prior to downloading the full dataset. If you use the dataset, please cite: "Phase-contrast X-ray tomography of free-breathing mice reveals the tracheobronchial tree", Kian Shaker, Ilian Häggmark, Jakob Reichmann, Marie Arsenian-Henriksson, and Hans M. Hertz, 2021 (under review)'. A red box highlights the 'Files' section, which contains a 'Restricted Access' heading, a paragraph explaining that access requests are subject to the owner's discretion, a text input field with the placeholder 'Dataset is available upon request. Please state the purpose of your request and the intended usage of the data.', and a 'Request access...' button.

zenodo Search Upload Communities

March 5, 2021 Dataset Restricted Access

Phase-contrast X-ray tomography of free-breathing murine lungs

Kian Shaker; Ilian Häggmark; Jakob Reichmann; Marie Arsenian-Henriksson; Hans M. Hertz

Full resolution phase-contrast X-ray tomographic dataset acquired of the lungs of a free-breathing mouse (NMRI nude mice, BomTac:NMRI-Foxn1tm, Taconic Biosciences, DK), weighing roughly 28 g.

Datasize: 3851x3951 pixels per slice, stack of 1700 slices, 16-bit, .tif

Voxel-size: 8.25x8.25x8.25 micrometer

A sample slice is available for download, prior to downloading the full dataset.

If you use the dataset, please cite:

"Phase-contrast X-ray tomography of free-breathing mice reveals the tracheobronchial tree", Kian Shaker, Ilian Häggmark, Jakob Reichmann, Marie Arsenian-Henriksson, and Hans M. Hertz, 2021 (under review)

Files

Restricted Access

You may request access to the files in this upload, provided that you fulfil the conditions below. The decision whether to grant/deny access is solely under the responsibility of the record owner.

Dataset is available upon request. Please state the purpose of your request and the intended usage of the data.

Request access...

Accessible (A1) - comments

FAIR data \neq open data!

- Common misconception stemming from A1.1
- FAIR data *can* be open data, but it has nothing to do with A1

Access can/should be realized over APIs

- SPARQL endpoints
- HTTP APIs
- Client libraries

Access is not only “click to download”

Tombstone pages (A2)

Metadata is accessible, even when the data is no longer available

The screenshot shows the Harvard Dataverse interface for a dataset titled "2000 Utah Colleges Exit Poll". The page includes a navigation bar with "Add Data", "Search", "About", "User Guide", "Support", "Sign Up", and "Log In". Below the navigation bar, the breadcrumb "Harvard Dataverse > 2000 Utah Colleges Exit Poll" is visible. A "Contact" button is located in the top right. The dataset title "2000 Utah Colleges Exit Poll" is displayed with a document icon and a red "Deaccessioned" label. A text box provides the citation: "David B. Magleby; Howard B. Christensen; Scott D. Grimshaw, 2019, '2000 Utah Colleges Exit Poll', <https://doi.org/10.7910/DVN/2Z9KDF>, Harvard Dataverse, V1, DEACCESSIONED VERSION, UNF:6:ME7YkGved9FxnBuA4Ytw== [fileUNF]". A red box highlights the "Deaccession Reason" section, which states: "User error. Do not use. Look under CSED and Utah Colleges Exit Poll". Below this, a "Versions" section is partially visible. At the bottom, a table lists the dataset version 1.0, its summary (repeating the deaccession reason), contributors (CSED CSED), and the publication date (Dec 30, 2019).

HARVARD
Dataverse

Add Data ▾ Search ▾ About User Guide Support Sign Up Log In

Harvard Dataverse > 2000 Utah Colleges Exit Poll

Contact

2000 Utah Colleges Exit Poll
Deaccessioned

David B. Magleby; Howard B. Christensen; Scott D. Grimshaw, 2019, "2000 Utah Colleges Exit Poll", <https://doi.org/10.7910/DVN/2Z9KDF>, Harvard Dataverse, V1, DEACCESSIONED VERSION, UNF:6:ME7YkGved9FxnBuA4Ytw== [fileUNF]

Deaccession Reason
User error. Do not use. Look under CSED and Utah Colleges Exit Poll

Versions

Dataset	Summary	Contributors	Published
1.0	Deaccessioned Reason: User error. Do not use. Look under CSED and Utah Colleges Exit Poll	CSED CSED	Dec 30, 2019

Interoperable

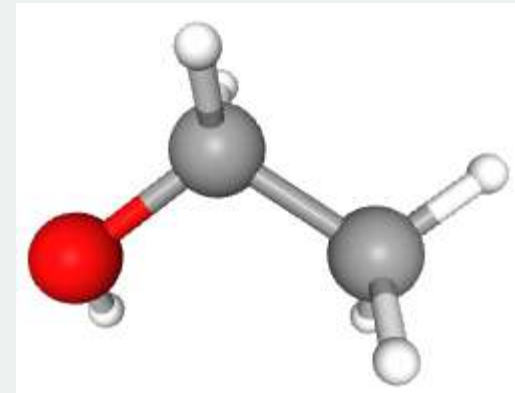
1. (Meta)data use a formal, accessible, shared, and broadly applicable **language for knowledge representation**.
2. (Meta)data use **vocabularies** that follow FAIR principles
3. (Meta)data include **qualified references** to other (meta)data

I1. (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation

“Data that should be readable for machines without the need for specialised or ad hoc algorithms, translators, or mappings”

Use:

- Common formats
 - RDF, JSON (+schema),
 - CSV (+ good README)
- Well defined/described data models
- Known representations
 - e.g. InChi Key: IDGUHHHCWSQLU-UHFFFAOYSA-N



I2: (Meta)data use **vocabularies** that follow the FAIR principles

Help evade ambiguities

“My plane lands in London...” – where exactly?

County	ICAO	IATA	Airport name	Usage
Greater London	EGKB	BQH	<i>London Biggin Hill Airport</i>	Public
Greater London	EGML		<i>Damyns Hall Aerodrome</i>	Private
Greater London	EGLL	LHR	<i>Heathrow Airport</i>	Public
Greater London	EGWU	NHT	<i>RAF Northolt</i>	Military
Greater London	EGLC	LCY	<i>London City Airport</i>	Public
Greater London	EGLW		<i>London Heliport</i>	Public

Controlled vocabularies: IATA and ICAO

https://en.wikipedia.org/wiki/List_of_airports_in_the_United_Kingdom_and_the_British_Crown_Dependencies

Vocabularies

- Tag units of information to make search and retrieval easier
- No need to be an expert in car manufacturers to rent a car worldwide
- Example: IDAD – Intermediate category, 4/5 doors, automatic transmission, diesel engine, air-conditioning fitted

Category	Type	Trans / Driven wheels	Fuel / air:con
M: Mini	B: 2-3 Door	M: Manual (drive unspecified)	R: Unspecified Fuel With Air
N: Mini Elite	C: 2/4 Door	N: Manual 4WD	N: Unspecified Fuel Without Air
E: Economy	D: 4-5 Door	C: Manual AWD	D: Diesel Air
H: Economy Elite	W: Wagon/Estate	A: Auto (drive unspecified)	Q: Diesel No Air
C: Compact	V: Passenger Van	B: Auto 4WD	H: Hybrid Air
D: Compact Elite	L: Limousine	D: Auto AWD	I: Hybrid No Air
I: Intermediate	S: Sport		E: Electric Air
J: Intermediate Elite	T: Convertible		C: Electric No Air
S: Standard	F: SUV		L: LPG/Compressed Gas Air
R: Standard Elite	J: Open Air All Terrain		S: LPG/Compressed Gas No Air
F: Fullsize	X: Special		A: Hydrogen Air
G: Fullsize Elite	P: Pick up Regular Cab		B: Hydrogen No Air
P: Premium	Q: Pick up Extended Cab		M: Multi Fuel/Power Air
U: Premium Elite	Z: Special Offer Car		F: Multi Fuel/Power No Air
L: Luxury	E: Coupe		V: Petrol Air
W: Luxury Elite	M: Monospace		Z: Petrol No Air
O: Oversize	R: Recreational Vehicle		U: Ethanol Air
X: Special	H: Motor Home		X: Ethanol No Air
	Y: 2 Wheel Vehicle		
	N: Roadster		
	G: Crossover		
	K: Commercial Van/Truck		

Vocabularies

Less time/money spent on data cleaning

- Different languages
- Spelling mistakes
- Abbreviations
- Capital letters

Vienna	Beč (Croatian, Serbian, older Bulgarian), Beç (older Turkish)*, Bech or Vidnya (Romani), Bécs (Hungarian)*, Bin / Pin - 빈 (Korean), Dunaj (Slovene)*, Fienna (Welsh), Vedunia (Celtic), Vena - Вена (Russian), Videň (Czech)*, Viden' / Videň (Ukrainian)*, Viedeň (Slovak), Viên (Vietnamese), Viena / Vijena/ Виена (Belarusian, Bulgarian, Macedonian), Viena (Catalan*, Lithuanian, Portuguese*, Romanian*, Spanish*, Tagalog*), Vienna (Italian)*, Vienne (French)*, Viénni - Βιέννη (Greek), Vieno (Esperanto), Viin (Estonian), Vin - ויין (Yiddish), Vín (Irish, Icelandic), Vina - וינה (Hebrew), Vínarborg (Icelandic variant), Vindobona (Latin), Vīne (Latvian)*, Viyana (Turkish)*, Vjenë (Albanian), Vjenna (Maltese), Vyana (Azeri), Wean (local Viennese, Austrian and Bavarian dialects)*, Weiyena - 維也納 (Chinese)*, Wene (Afrikaans), Wenen (Dutch)*, Wiedeń (Polish)*, Wien (Danish*, Finnish*, German*, Norwegian*, Swedish*), Wīn - ウィーン (Japanese)*, Wina (Indonesian), فيينا (Arabic), وين (Persian)
--------	--

Vocabularies (I2)

UniProtKB - O00559 (RCAS1_HUMAN)

Display [Help video](#) [BLAST](#) [Align](#) [Format](#) [Add to basket](#) [History](#) [Add a publication](#) [Feedback](#)

Entry **Protein** Receptor-binding cancer antigen expressed on SISO cells
Gene **EBAG9**
Organism *Homo sapiens (Human)*
Status [Reviewed](#) - Annotation score: [★★★★★](#) - Experimental evidence at protein level²

Function¹
May participate in suppression of cell proliferation and induces apoptotic cell death through activation of interleukin-1-beta converting enzyme (ICE)-like proteases.
[3 Publications](#)

Miscellaneous
May serve as a prognostic marker for cancers such as adenocarcinomas of the lung and breast cancers. It is present and overexpressed in many patients suffering from breast carcinomas, its level of expression correlates with tumor grade, suggesting that it may be involved in cancer immune escape. According to PubMed:12672804, it is however not directly a tumor-associated antigen, but it rather modulates surface expression of tumor-associated O-linked glycan Tn when it is overexpressed, suggesting that it contributes indirectly to the antigenicity of tumor cells.

Caution
It was initially reported to be a ligand for some putative receptor present on T-, B-, natural killer (NK) cells and various human cell lines. However, PubMed:12672804 showed that it does not bind any receptor. [Curated](#)

GO - Molecular function¹
• peptidase activator activity involved in apoptotic process [Source: UniProtKB](#)

Complete GO annotation on QuickGO ...

GO - Biological process¹
• regulation of cell growth [Source: UniProtKB](#)

Complete GO annotation on QuickGO ...

Keywords¹
Biological process: [Apoptosis](#)

Enzyme and pathway databases
PathwayCommons¹: O00559
Reactome¹: R-HSA-9018519, Estrogen-dependent gene expression

Names & Taxonomy¹

Protein names ¹	Recommended name: Receptor-binding cancer antigen expressed on SISO cells Alternative name(s): • Cancer-associated surface antigen RCAS1 • Estrogen receptor-binding fragment-associated gene 9 protein
Gene names ¹	Name: EBAG9 Synonyms: RCAS1
Organism ¹	<i>Homo sapiens (Human)</i>
Taxonomic identifier ¹	9606 [NCBI]
Taxonomic lineage ¹	Eukaryota › Metazoa › Chordata › Craniata › Vertebrata › Euteleostomi › Mammalia › Eutheria › Euarchontoglires › Primates › Haplorhini › Catarrhini › Hominoidea › Homo
Proteomes ¹	UP00005640 Component 1: Chromosome 8

Vocabularies (I2)

Each metadata field has its definition

Organism

Last modified April 10, 2018

This subsection of the [Names and taxonomy](#) section provides information on the name(s) of the organism that is the source of the protein sequence.

The organism designation consists of the Latin scientific name, usually composed of the genus and species names (the binomial system developed by Linnaeus), followed optionally by the English common name and a synonym.

Examples: *Bacillus subtilis*, *Homo sapiens* (Human), *Cardamine pratensis* (Cuckoo flower) (Alpine bitter cross)

The synonym can be a common name in English (or in Latin in the case of some historical legacy names).

Example: *Radianthus magnifica* (Magnificent sea anemone) (*Heteractis magnifica*)

In the case of viruses, the designation does not follow the binomial system. The English common name is used as the scientific name, sometimes followed by an acronym. When possible, viruses are named according to the nomenclature of the International Committee on Taxonomy of Viruses (ICTV).

Examples: Human immunodeficiency virus type 1 (isolate BRU/LAI group M subtype B) (HIV-1), Influenza A virus (strain A/Aichi/2/1968 H3N2)

The organism name can differ from that given by the international nucleotide sequence databases for the same taxon. This is mainly due to our efforts in providing the most descriptive common names and synonyms to our users.

Note that the proteome for a given organism, when available, can be accessed through the [proteomes](#) page of our website.

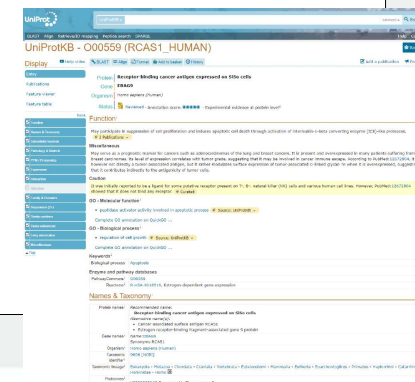
Related documents

[Taxonomy](#)

[Controlled vocabulary of species](#)

[What are proteomes?](#)

[What are reference proteomes?](#)



Vocabularies (I2)

Taxonomy - Homo sapiens (Human) (SPECIES)

Map to UniProtKB (194,609)
Reviewed (20,397)
Swiss-Prot
Unreviewed (174,212)
TrEMBL
Proteomes (3)

Format

Mnemonic	HUMAN
Taxon identifier	9606
Scientific name	Homo sapiens
Taxonomy navigation	↑ Homo ↓ Choose one All lower taxonomy nodes (2)
Common name	Human
Synonym	-
Other names	> Home sapiens > Homo samplens > Homo sapeins > Homo saplan > Homo sapians More >
Rank	SPECIES
Lineage	> cellular organisms > Eukaryota > Opisthokonta > Metazoa > Eumetazoa > Bilateria > Deuterostomia > Chordata > Craniata > Vertebrata > Gnathostomata > Teleostomi > Euteleostomi > Sarcopterygii > Dipnotetrapodomorpha > Tetrapoda > Amniota > Mammalia > Theria > Eutheria > Boreoeutheria > Euarchontoglires > Primates > Haplorhini > Simiiformes > Catarrhini > Hominoidea > Hominidae > Hominiinae > Homo

Each metadata **value** comes from a controlled vocabulary – no free form answers.

UniProt - UniProtKB - O00559 (RCAS1_HUMAN)

Display: [View](#) [Print](#) [Download](#) [Share](#) [Feedback](#) [Help](#)

Accession: **O00559** (RCAS1_HUMAN)

Protein name: **Receptor tyrosine kinase expressed on S16 cells**

Gene name: **RCAS1**

Organism: **Homo sapiens (Human)**

Keywords: **Receptor** **Autophosphorylation** **Experimental evidence of protein function**

Function: **May participate in suppression of cell proliferation and telomere apoptosis, cell death through activation of telomerase-inhibiting complex (TIN2-like protein).**

Biological process: **May serve as a protein kinase for cancer, such as glioblastoma of the lung and breast cancer. It is present and overexpressed in many cell lines suffering from these cancerous cells that is dependent on telomerase activity. Suggests that it may be involved in several cancer events, including in telomerase (TIN2). It is thought to directly or indirectly interact with telomerase, suggesting that it may be involved in several cancer events, including in telomerase (TIN2). It is thought to directly or indirectly interact with telomerase, suggesting that it may be involved in several cancer events, including in telomerase (TIN2). It is thought to directly or indirectly interact with telomerase, suggesting that it may be involved in several cancer events, including in telomerase (TIN2).**

Cellular component: **Cell membrane**

Substrate: **ATP**

Enzyme class: **EC: 2.7.10.1**

Enzyme name: **EC: 2.7.10.1**

Enzyme reaction: **ATP + H₂O → ADP + Pi**

Enzyme regulation: **Complete cell activation on Q00559 ...**

Keywords: **Receptor**

Biological process: **Regulation**

Enzyme and pathway databases: **UniProtKB: O00559**

Names & Taxonomy

Protein name: **Receptor tyrosine kinase expressed on S16 cells**

Gene name: **RCAS1**

Organism: **Homo sapiens (Human)**

Enzyme class: **EC: 2.7.10.1**

Enzyme name: **EC: 2.7.10.1**

Enzyme reaction: **ATP + H₂O → ADP + Pi**

Enzyme regulation: **Complete cell activation on Q00559 ...**

Keywords: **Receptor**

Biological process: **Regulation**

Enzyme and pathway databases: **UniProtKB: O00559**

Names & Taxonomy

Protein name: **Receptor tyrosine kinase expressed on S16 cells**

Gene name: **RCAS1**

Organism: **Homo sapiens (Human)**

Enzyme class: **EC: 2.7.10.1**

Enzyme name: **EC: 2.7.10.1**

Enzyme reaction: **ATP + H₂O → ADP + Pi**

Enzyme regulation: **Complete cell activation on Q00559 ...**

Keywords: **Receptor**

Biological process: **Regulation**

Enzyme and pathway databases: **UniProtKB: O00559**

Qualified References (I3)

Meaningful links to describe connections

- Dataset X was *derived from* dataset Y
- Dataset Y was *produced* using code Z

Standard relations define by Data Cite

- <https://schema.datacite.org/meta/kernel-4.3/>

Use persistent identifiers

```
<language>en</language>
<resourceType resourceTypeGeneral="Workflow">Software</resourceType>
<relatedIdentifiers>
  <relatedIdentifier relatedIdentifierType="DOI"
    relationType="IsReferencedBy">10.5072/2047-217X-1-1</relatedIdentifier>
  <relatedIdentifier relatedIdentifierType="DOI"
    relationType="Compiles">10.5072/100038</relatedIdentifier>
</relatedIdentifiers>
<sizes>
  <size>31 MB</size>
</sizes>
```

IsContinuedBy
Continues
IsDescribedBy
Describes
HasMetadata
IsMetadataFor
HasVersion
IsVersionOf
IsNewVersionOf
IsPreviousVersionOf
IsPartOf
HasPart
IsReferencedBy
References
IsDocumentedBy
Documents
IsCompiledBy
Compiles
IsVariantFormOf
IsOriginalFormOf
IsIdenticalTo
IsReviewedBy
Reviews
IsDerivedFrom
IsSourceOf
IsRequiredBy
Requires
IsObsoletedBy
Obsoletes

Qualified References (I3)

The screenshot shows a dataset page on the TU Wien website. The title is "European Sentinel-1 Forest Type and Tree Cover Density Maps". The authors listed are Dostalova, Alena; Cao, Senmao; and Wagner, Wolfgang. The page includes a description of the dataset, a dataset record, code availability information, and acknowledgements. A "Details" sidebar on the right lists the resource type as "Dataset", formats as "application/x-geotiff", and provides several related identifiers (DOI, Zenodo, GitHub, and references).

January 19, 2021 | Version 1.0

European Sentinel-1 Forest Type and Tree Cover Density Maps

Dostalova, Alena¹; Cao, Senmao^{1,2}; Wagner, Wolfgang^{1,2} [show affiliations](#)

Description
This dataset was generated by the TU Wien Department of Geodesy and Surveying.

European Sentinel-1 forest type and tree cover density maps represent the dominant forest type class (conifer, deciduous, or mixed) and the percentage of forest canopy cover within the 100 m pixel. The forest type map shows the dominant forest type class (conifer, deciduous, or mixed) and the percentage of forest canopy cover within the 100 m pixel.

Please be referred to our peer-reviewed article at <https://doi.org/10.1016/j.cageo.2014.07.005> for an assessment across Europe.

Dataset Record

The forest type and tree cover density maps are sampled at 10 m, georeferenced to the Equi7Grid and divided into square tiles of 10 m. The maps consist of 728 tiles over the European continent, with data v

The tiles' file-format is a LZW-compressed GeoTIFF holding 16-bit and georeference. Compatibility with common geographic information libraries as GDAL is given.

In this repository, we provide each forest map as tiles, whereas two download below.

Code Availability

For the usage of the **Equi7Grid** we provide data and tools via the python package available on GitHub at <https://github.com/TUW-GEO/Equi7Grid>. More details on the grid reference can be found in <https://www.sciencedirect.com/science/article/pii/S0098300414001629>.

Acknowledgements

The computational results presented have been achieved using the Vienna Scientific Cluster (VSC).

Details

Licenses

Resource type Dataset

Formats application/x-geotiff

Related identifiers

isreferencedby	10.3390/rs13030337 (doi)
issupplementto	10.5281/zenodo.3515933 (doi) https://github.com/TUW-GEO/Equi7Grid (url)
references	10.1080/01431161.2018.1479788 (doi) 10.1016/j.cageo.2014.07.005 (doi)

Paper citing this dataset

Code

Paper describing the method to produce this dataset

<https://researchdata.dl.hpc.tuwien.ac.at/records/tkkfs-11b75>

Qualified References (I3)

README.rst

🔗 Equi7Grid

build passing coverage 32% pypi package 0.0.12 docs passing

A python class for working with Equi7Grid - how to convert to - how to use the tiling system - etc.

It's a python package that handles the geometric and geographic operations of a gridded and tiled projection system. It was designed for data cubes ingesting satellite imagery and builds the basis for the Equi7Grid (see <https://github.com/TUW-GEO/Equi7Grid>).

A detailed documentation on the Equi7Grid definition is at:

`~/docs/doc_files/`

Overlays for visualisation in Google Earth can be found here:

`~/docs/doc_files/google_earth_overlays/`

Citation

DOI [10.5281/zenodo.1048530](https://doi.org/10.5281/zenodo.1048530)

If you use the software in a publication then please cite it using the Zenodo DOI. Be aware that this badge links to the latest package version.

Please select your specific version at <https://doi.org/10.5281/zenodo.1048530> to get the DOI of that version. You should normally always use the DOI for the specific version of your record in citations. This is to ensure that other researchers can access the exact research artefact you used for reproducibility.

You can find additional information regarding DOI versioning at <http://help.zenodo.org/#versioning>

Qualified References (I3)

The image shows a screenshot of the UniProtKB entry for O00559 (RCAS1_HUMAN). The page is divided into several sections. A red box highlights the 'Cross-references' section, which is currently empty. Another red box highlights the entry 'AAB61617.1' in the 'Sequence databases' section, which is linked to the mRNA translation of AF006265. The 'Function' section describes the protein's role in cell proliferation and apoptosis. The 'GO - Molecular function' section lists 'peptidase activator activity involved in apoptotic process'. The 'GO - Biological process' section lists 'regulation of cell growth'. The 'Names & Taxonomy' section provides details on the protein's name, gene name (EBAG9), and taxonomic classification (Homo sapiens).

UniProtKB - O00559 (RCAS1_HUMAN)

Cross-referencesⁱ

Web resourcesⁱ

[Atlas of Genetics and Cytogenetics in Oncology and Haematology](#)

Sequence databases

Select the link destinations:

- EMBLⁱ
- GenBankⁱ
- ODBJⁱ

AF006265 mRNA Translation: **AAB61617.1**

AB007619 mRNA Translation: BAA22572.1

AY653072 mRNA Translation: AAU85838.1

AK290651 mRNA Translation: BAF83340.1

CR456984 mRNA Translation: CAG33265.1

AC079061 Genomic DNA No translation available.

AP000427 Genomic DNA No translation available.

BC005249 mRNA Translation: AAH05249.1

BC017729 mRNA Translation: AAH17729.1

BC022506 mRNA Translation: AAH22506.1

CCDSⁱ CCDS6313.1 [O00559-1]

RefSeqⁱ NP_001265867.1, NM_001278938.1 [O00559-1]
NP_004206.1, NM_004215.4 [O00559-1]
NP_936056.1, NM_198120.2 [O00559-1]
XP_016869449.1, XM_017013960.1 [O00559-1]

3D structure databases

ModBaseⁱ [Search...](#)

SWISS-MODEL-Workspaceⁱ [Submit a new modelling project...](#)

Protein-protein interaction databases

BioGRIDⁱ 114607, 40 interactors

IntActⁱ O00559, 30 interactors

STRINGⁱ 9606.ENSPP00000337675

Reusable

R1. (Meta)data are richly described with a plurality of accurate and relevant attributes

- R1.1. (Meta)data are released with a clear and accessible data usage **license**
- R1.2. (Meta)data are associated with detailed **provenance**
- R1.3. (Meta)data meet domain-relevant **community standards**

R1.1. (Meta)data are released with a clear and accessible data usage license

Public repository on GitHub

- May suggest that authors are willing to share code

No license

- no possibility for reuse
- can only be viewed (only because terms of use enforce that)

Code without a license is like an object in a museum

- You can watch and admire it, but you cannot touch it!

More on licenses in other lectures

License (R1.1)

The screenshot displays the GitHub repository page for PyTorch. At the top, the repository name 'pytorch / pytorch' is shown along with statistics: 23.7k users, 1.4k watches, 36.6k stars, and 9.3k forks. Below this, navigation tabs include Code, Issues (3,951), Pull requests (1,284), Actions, Projects (5), Wiki, Security, and Insights. The repository description is 'Tensors and Dynamic neural networks in Python with strong GPU acceleration' with a link to 'https://pytorch.org'. A list of tags includes 'neural-network', 'autograd', 'gpu', 'numpy', 'deep-learning', 'tensor', 'python', and 'machine-learning'. A summary bar shows 24,578 commits, 2,993 branches, 0 packages, 31 releases, and 1,316 contributors. A red box highlights the 'View license' button in this bar. Below the summary, there are buttons for 'Branch: master', 'New pull request', 'Create new file', 'Upload files', 'Find file', and 'Clone or download'. A commit history table is visible, with a red box highlighting the commit 'LICENSE' by 'jamesr66a and facebook-github-bot' with the message 'Move copyright lines back to NOTICE file, fixes #6911 (#8310)' from 2 years ago. At the bottom, a 'License' section is highlighted with a red box, containing the text: 'PyTorch is BSD-style licensed, as found in the LICENSE file.'

Search or jump to... Pull requests Issues Marketplace Explore

pytorch / pytorch Used by 23.7k Watch 1.4k Star 36.6k Fork 9.3k

Code Issues 3,951 Pull requests 1,284 Actions Projects 5 Wiki Security Insights

Tensors and Dynamic neural networks in Python with strong GPU acceleration <https://pytorch.org>

neural-network autograd gpu numpy deep-learning tensor python machine-learning

24,578 commits 2,993 branches 0 packages 31 releases 1,316 contributors [View license](#)

Branch: master New pull request Create new file Upload files Find file Clone or download

	jamesr66a and facebook-github-bot [quantization] Make FP16 RNN use new prepack op (#34339) ...	Latest commit 8a17dc6 7 minutes ago
	.circleci: Remove macOS builds related to CUDA (#34333)	21 hours ago
	.ctags.d: Add a .ctags.d/ toplevel directory (#18827)	11 months ago
	LICENSE: Move copyright lines back to NOTICE file, fixes #6911 (#8310)	2 years ago

License

PyTorch is BSD-style licensed, as found in the LICENSE file.

R1.2 (Meta)data are associated with detailed provenance

Provenance

- Describes origin of data
- Who? What? When? How?

Supports evaluation and can build trust in data

- ‘Officially, North Korea claims to have identified zero cases of COVID-19 inside its territory’

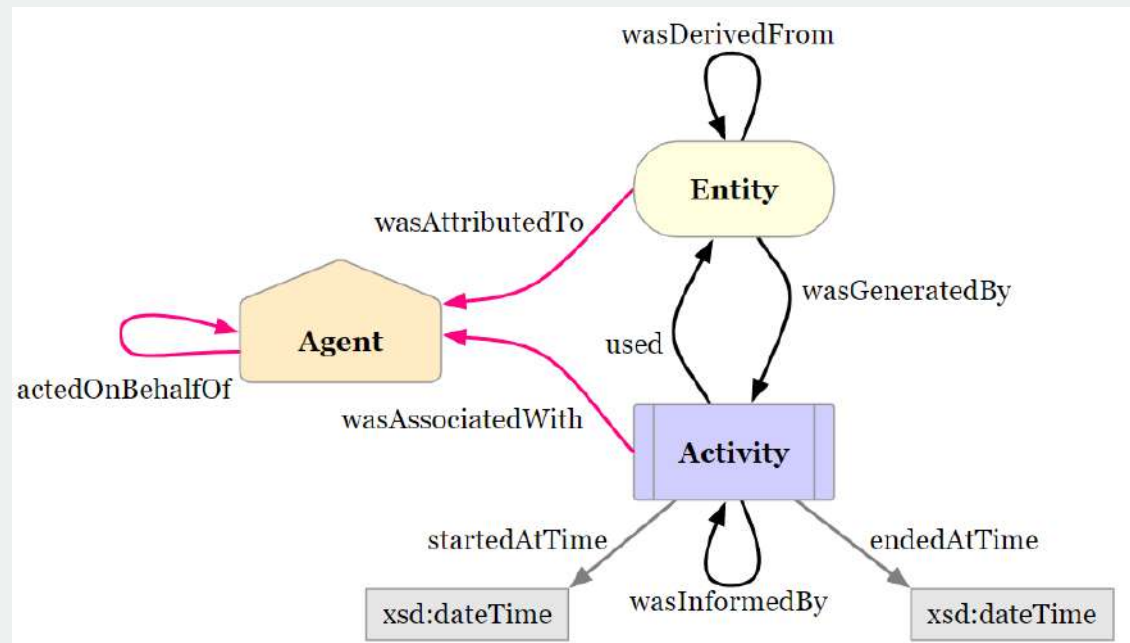
<https://www.npr.org/sections/goatsandsoda/2020/02/20/807027901/north-korea-claims-zero-cases-of-coronavirus-infection-but-experts-are-skeptical?t=1615196582563>



R1.2 (Meta)data are associated with detailed provenance

PROV-O: The PROV Ontology

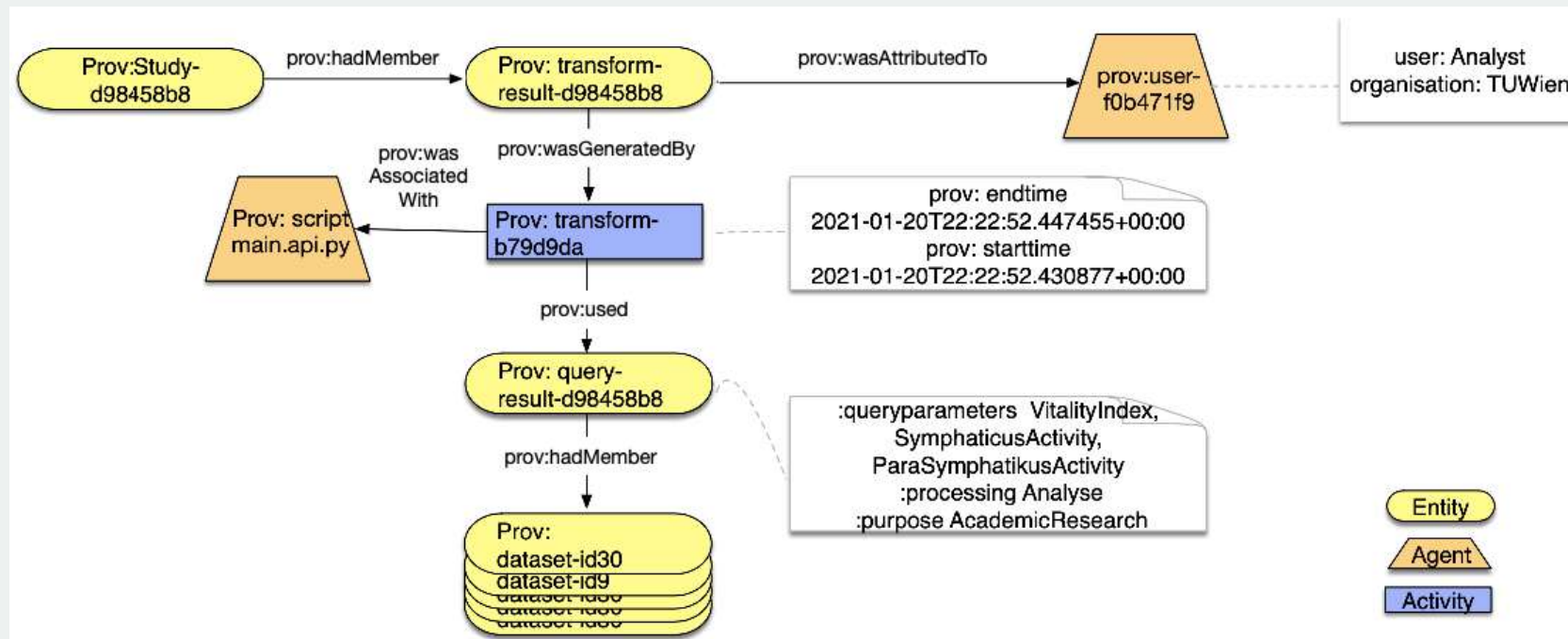
- Machine-actionable way to express provenance



R1.2 (Meta)data are associated with detailed provenance

Example of PROV-O Instance

- **Analyst** from TU Wien made a **study** in which data was **transformed**.
- To do so, a **script** (...) was executed to **query** data, using **parameters** (...) and following **datasets** were used (...).



R1.3. (Meta)data meet domain-relevant community standards

Who is the “community” ?

What is the “standard” ?

- English vs other languages

Metadata

- Domain independent
 - e.g. Dublin Core
- Domain specific
 - e.g. EXIF for images

Sometimes no common standard exist

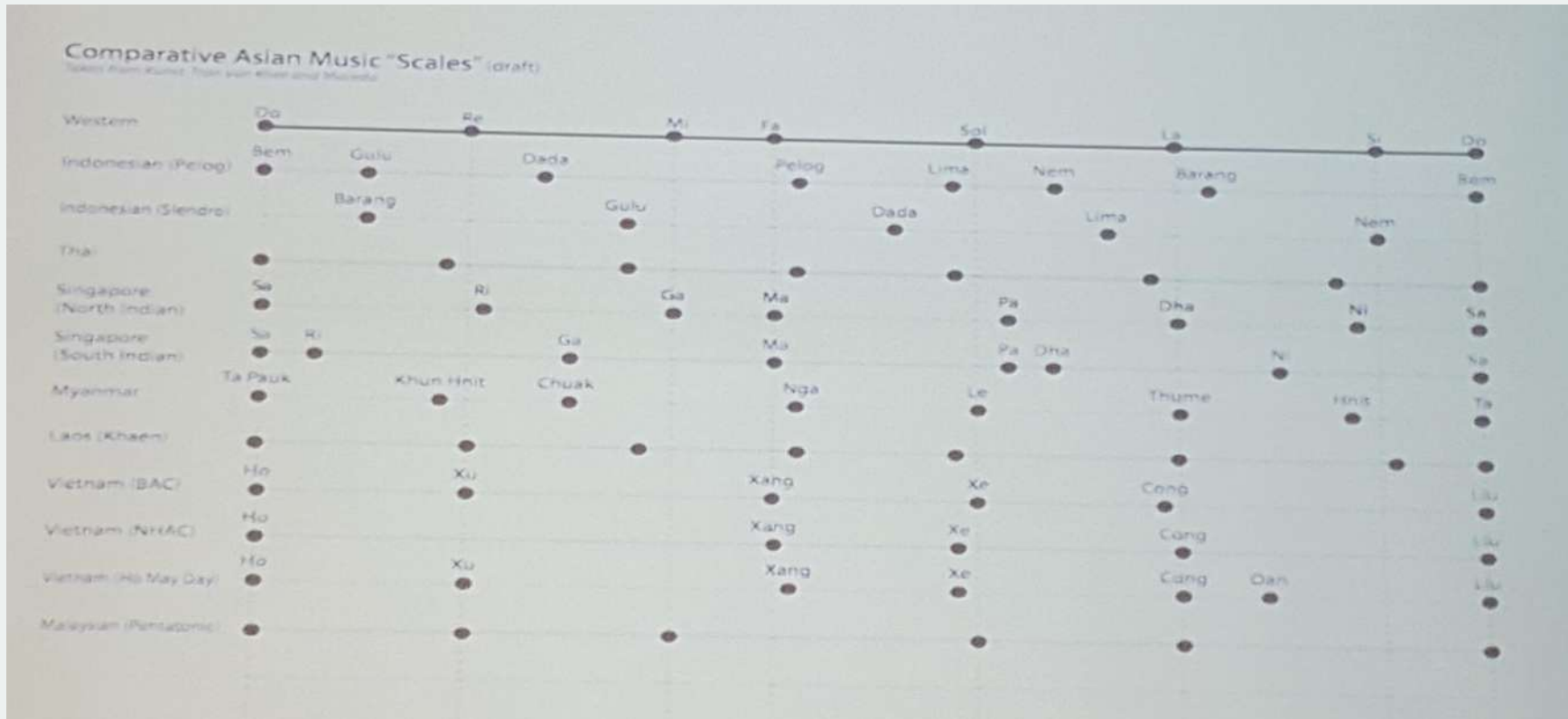
- Good documentation and README

There is no universal guideline – it always depends!

R1.3. (Meta)data meet domain-relevant community standards

Do Re Mi Fa Sol La Si Do

- Does not have to be a standard for everyone!



R1.3. (Meta)data meet domain-relevant community standards

Follow standards and domain specific conventions

Examples

- Sharing COBOL code
 - with Data Science students?
 - with mainframe operators?
- Obfuscated code/data

```
000024
000025 PROCEDURE DIVISION.
000026 0001-MAIN.
000027     INSPECT FUNCTION REVERSE (STR-1)
000028           TALLYING WS-LEN1 FOR LEADING SPACES.
000029     COMPUTE WS-LEN = LENGTH OF STR-1 - WS-LEN1.
000030     DISPLAY WS-LEN.
000031     MOVE 1 TO I.
000032     MOVE WS-LEN TO J.
000033     PERFORM REV-PARA WS-LEN TIMES.
000034     DISPLAY STR-1.
000035     DISPLAY STR-2.
000036     GOBACK.
000037     REV-PARA.
000038     MOVE STR-1(J:1) TO STR-2(I:1).
000039     SUBTRACT 1 FROM J.
000040     ADD 1 TO I.
000041     EXIT.
***** ***** Bottom of Data *****
```



R1.3. (Meta)data meet domain-relevant community standards

Good documentation supports reuse

- Removes ambiguities (especially where there are no common controlled vocabularies or others standards)

Example

- Confirmed cases of COVID-19: testing date vs reporting date

Indicators	Definition
Tests	Cumulative number of tests carried out for SARS-CoV-2, from 27 February 2020 up to and including the reporting date. Responsible for data consolidation: Office of the respective federal state government (Land), data status: morning of the reporting day.
Laboratory-confirmed cases	Cumulative number of laboratory-confirmed cases of SARS-CoV-2 infection (sum of "Active cases", "Recovered cases" and "Deceased cases") with laboratory diagnosis date since 27.02.2020 up to and including the reporting date .
Active cases	Cumulative number of laboratory-confirmed cases of SARS-CoV-2 infection with laboratory diagnosis date from 27.02.2020 up to and including the reporting date, which have not been classified as "recovered" or "deceased" on the reporting date.
Recovered cases	Cumulative number of laboratory-confirmed cases of SARS-CoV-2 infection with laboratory diagnosis date from 27.02.2020 up to and including the report date, which are classified as "recovered" on the report date. Definition of "recovered" (since 9 July): in the case of home care, 10-day home isolation after the onset of symptoms or laboratory diagnosis; in case of severe disease progression, the earliest 10 days after onset of symptoms, at least 48 hours without symptoms AND the following result by RT-PCR according to the Charité protocol: no nucleic acid detection of beta-coronavirus SARS-CoV-2 or nucleic acid detection of beta-coronavirus SARS-CoV-2 at a Ct value of more than 30. Further details can be found in the recommendation for the release of COVID-19 cases, recommendation for the release of COVID-19 cases from isolation.
Deceased cases	Cumulative number of laboratory-confirmed cases of SARS-CoV-2 infection with a laboratory diagnosis date from 27.02.2020 up to and including the report date, which are classified as "deceased" on the report date. Definition of "deceased": COVID-19 death is defined, for surveillance purposes, as one laboratory-confirmed case of COVID-19 resulting

FAIR ASSESSMENT

FAIR Assessment

Still pretty much work in progress

Most approaches

- Try to quantify FAIRness
- Based on self-assessment
 - <https://ardc.edu.au/resources/working-with-data/fair-data/fair-self-assessment-tool/>
 - <https://fairaware.dans.knaw.nl>



FAIR Aware – self-assessment example

F
A
I
R

FAIR questions

FINDABLE

1. Are you aware that a dataset should be assigned a globally unique persistent and resolvable identifier when deposited with a data repository? Yes No
2. Are you aware that when you deposit a dataset with a repository, you will need to provide some details (known as discovery metadata) in order to make the data findable, understandable and reusable to others? Yes No
3. Are you aware that the repository providing access to your dataset should make the metadata describing your datasets available in a format readable by machines as well as humans? Yes No

ACCESSIBLE

4. Are you aware that access to your dataset may need to be controlled and that metadata should include licence information under which the data can be reused? Yes No
5. Are you aware that metadata should remain available over time, even if the data is no longer accessible? Yes No

INTEROPERABLE

6. Are you aware that the metadata describing your datasets should use semantic vocabularies? Yes No

REUSABLE

7. Are you aware that provenance information about the collection and/or generation of data should be included in the metadata? Yes No
8. Are you aware that metadata describing your data should follow the specifications of a community-endorsed standard? Yes No
9. Are you aware that data should be deposited preferably in a file format that is open – to support reuse – and supported by the repository for long-term preservation? Yes No
10. Are you aware that maintaining your dataset FAIR over time requires professional data curation and preservation? Yes No

Findable

- Does the dataset have any identifiers assigned?
- Is the dataset identifier included in all metadata records/files describing the data?
- How is the data described with metadata?
- What type of repository or registry is the metadata record in?

Accessible

- How accessible is the data?
- Is the data available online without requiring specialised protocols or tools once access has been approved?
- Will the metadata record be available even if the data is no longer available?

Interoperable

- What (file) format(s) is the data available in?
- What best describes the types of vocabularies/ontologies/tagging schemas used to define the data elements?
- How is the metadata linked to other data and metadata (to enhance context and clearly indicate relationships)?

Reusable

- Which of the following best describes the licence/usage rights attached to the data?
- How much provenance information has been captured to facilitate data reuse?

RDA FAIR Data Maturity Model Specification and Guidelines Recommendation

17 minimum viable metrics to systematically measure the extent to which **research data objects** are FAIR

Not systems as a whole!

2.1 Globally Unique Identifier

FIELD	DESCRIPTION
Metric Identifier	FsF-F1-01D
Metric Name	Data is assigned a globally unique identifier.
Description	A data object may be assigned with a globally unique identifier such that it can be referenced unambiguously by humans or machines. Globally unique means an identifier should be associated with only one resource at any time. Examples of unique identifiers of data are Internationalized Resource Identifier (IRI) ³⁶ , Uniform Resource Identifier (URI) such as URL and URN, Digital Object Identifier (DOI), the Handle System, identifiers.org, w3id.org and Archival Resource Key (ARK). A data repository may assign a globally unique identifier to your data or metadata when you publish and make it available through its curation service.
FAIR Principle	F1. (Meta) data are assigned globally unique and persistent identifiers
CoreTrustSeal Alignment	R13. The repository enables users to discover the data and refer to them in a persistent way through proper citation
ASSESSMENT	
Requirement(s)	<ul style="list-style-type: none">• Data identifier (IRI, URL)• List of globally unique identifier schemes
Method	Check if the identifier is specified based on a globally unique identifier scheme.
COMMENTS	
Related Resources: <ul style="list-style-type: none">• Identifiers compiled by FAIRsharing, https://fairsharing.org/standards/?q=&selected_facets=type_exact:identifier%20schema• A list of Uniform Resource Identifier (URI) schemes, available in different formats, https://www.iana.org/assignments/uri-schemes/uri-schemes.xhtml#uri-schemes-1• Uniform Resource Identifier (URI) Generic Syntax (RFC 3986), https://tools.ietf.org/html/rfc3986	



<http://doi.org/10.5281/zenodo.4081213> (Published: 12 October 2020)



F-UJI

Automated FAIR Data Assessment Tool

Disclaimer: The test results shown here are based on preliminary data and code which still is under development. F-UJI is rapidly evolving and not yet available in a productive environment.

Research Data Object (URL/PID):*

OAI-PMH:

Enable caching? ⓘ Use DataCite? ⓘ

Assessment Results:

Evaluated Resource:

Daily snow cover grid maps over Austria in the period 2000-2020

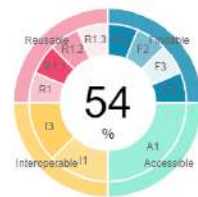
Resource PID/URL: <https://researchdata.tuwien.ac.at/records/4wmcs-ed919>

Metric Version: metrics_v0.4

Metric Specification: <https://doi.org/10.5281/zenodo.4081213>

Software version: v1.3.5b

Summary:



Findable: 5 of 7



Accessible: 1 of 3



Interoperable: 3 of 4



Reusable: 4 of 10



Code still under development!

<https://www.f-uji.net/index.php>

Report:

Findable

FsF-F1-01D - Data is assigned a globally unique identifier. 

FsF-F1-02D - Data is assigned a persistent identifier. 

FsF-F2-01M - Metadata includes descriptive core elements (creator, title, data identifier, publisher, publication date, summary and keywords) to support data findability. 

FsF-F3-01M - Metadata includes the identifier of the data it describes. 

FsF-F4-01M - Metadata is offered in such a way that it can be retrieved programmatically. 

Accessible

FsF-A1-01M - Metadata contains access level and access conditions of the data. 

FsF-A1-03D - Data is accessible through a standardized communication protocol. 

FsF-A1-02M - Metadata is accessible through a standardized communication protocol. 

Interoperable

FsF-I1-01M - Metadata is represented using a formal knowledge representation language. 

FsF-I1-02M - Metadata uses semantic resources. 

FsF-I3-01M - Metadata includes links between the data and its related entities. 

Reusable

FsF-R1-01MD - Metadata specifies the content of the data. 









FsF-R1.1-01M - Metadata includes license information under which data can be reused. 

FsF-R1.2-01M - Metadata includes provenance information about data creation or generation. 

FsF-R1.3-01M - Metadata follows a standard recommended by the target research community of the data. 

FsF-R1.3-02D - Data is available in a file format recommended by the target research community. 

Metric tests:

Test:	Test name:	Result:
FsF-R1-01MD-1	Minimal information about available data content is given in metadata	
FsF-R1-01MD-1a	Resource type (e.g. dataset) is given in metadata	
FsF-R1-01MD-1b	Information about data content (e.g. links) is given in metadata	
FsF-R1-01MD-2	Verifiable data descriptors (file info, measured variables or observation types) are specified in metadata	
FsF-R1-01MD-2a	File size and type information are specified in metadata	
FsF-R1-01MD-2b	Measured variables or observation types are specified in metadata	
FsF-R1-01MD-3	Data content matches file type and size specified in metadata	
FsF-R1-01MD-4	Data content matches measured variables or observation types specified in metadata	

Debug:

Level:	Message:
INFO	Object landing page accessible status -: True
SUCCESS	Resource type specified -: dataset
WARNING	NO data object content available/accessible to perform file descriptors (type and size) tests
WARNING	NO measured variables found in metadata, skip 'measured_variable' test.
WARNING	Measured variables given in metadata do not match data object content

FAIR Digital Object

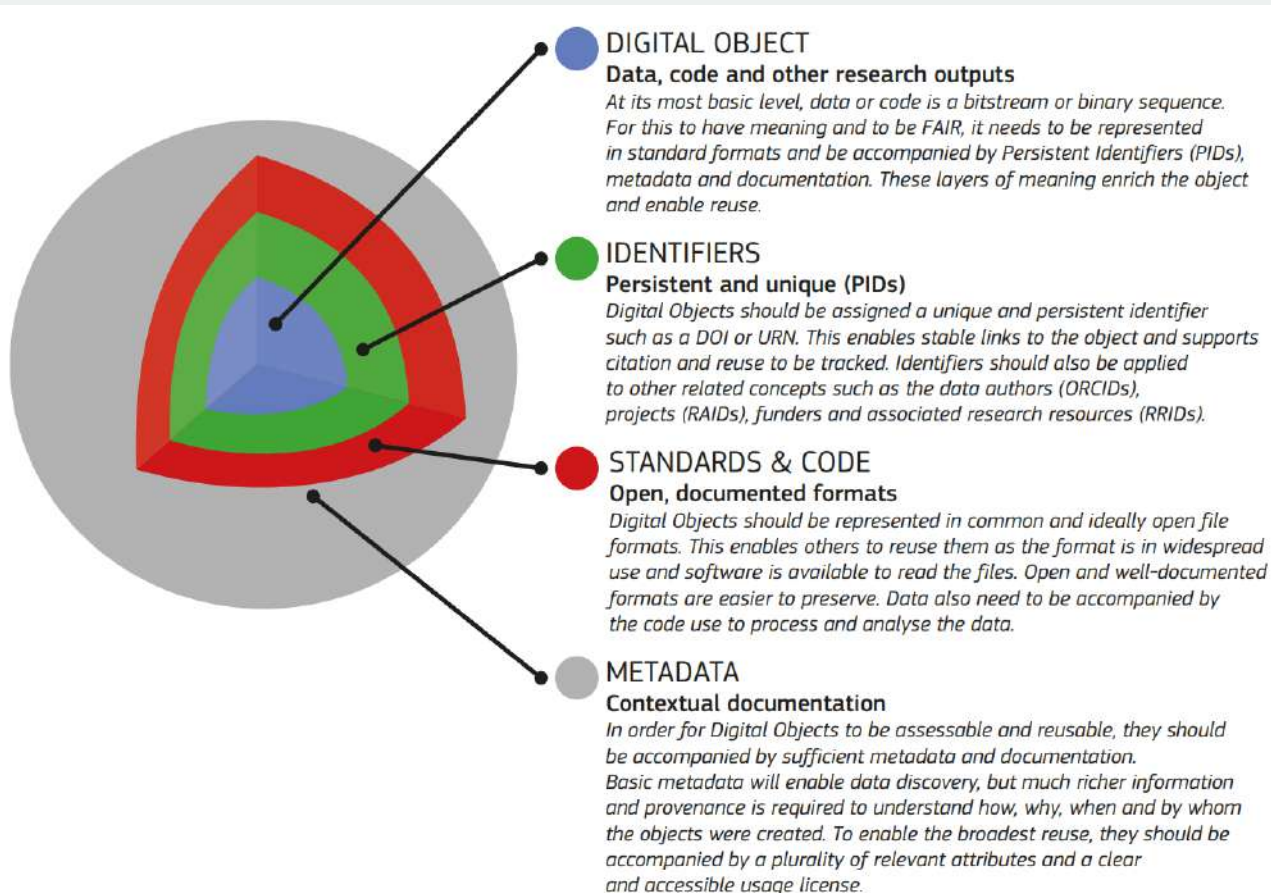


Figure 8. A model for FAIR Digital Objects, noting the elements that need to be in place for data to be Findable, Accessible, Interoperable and Reusable

Turning FAIR into reality <https://op.europa.eu/s/oM5N>



RO-Crate as an example of a FAIR Data Object

Currently, no commonly accepted FAIR Data Object

RO-Crate

- packages **research artefacts** along with their **metadata** in a **machine readable** way
- based on Schema.org annotations in JSON-LD
- <https://www.researchobject.org/ro-crate/>

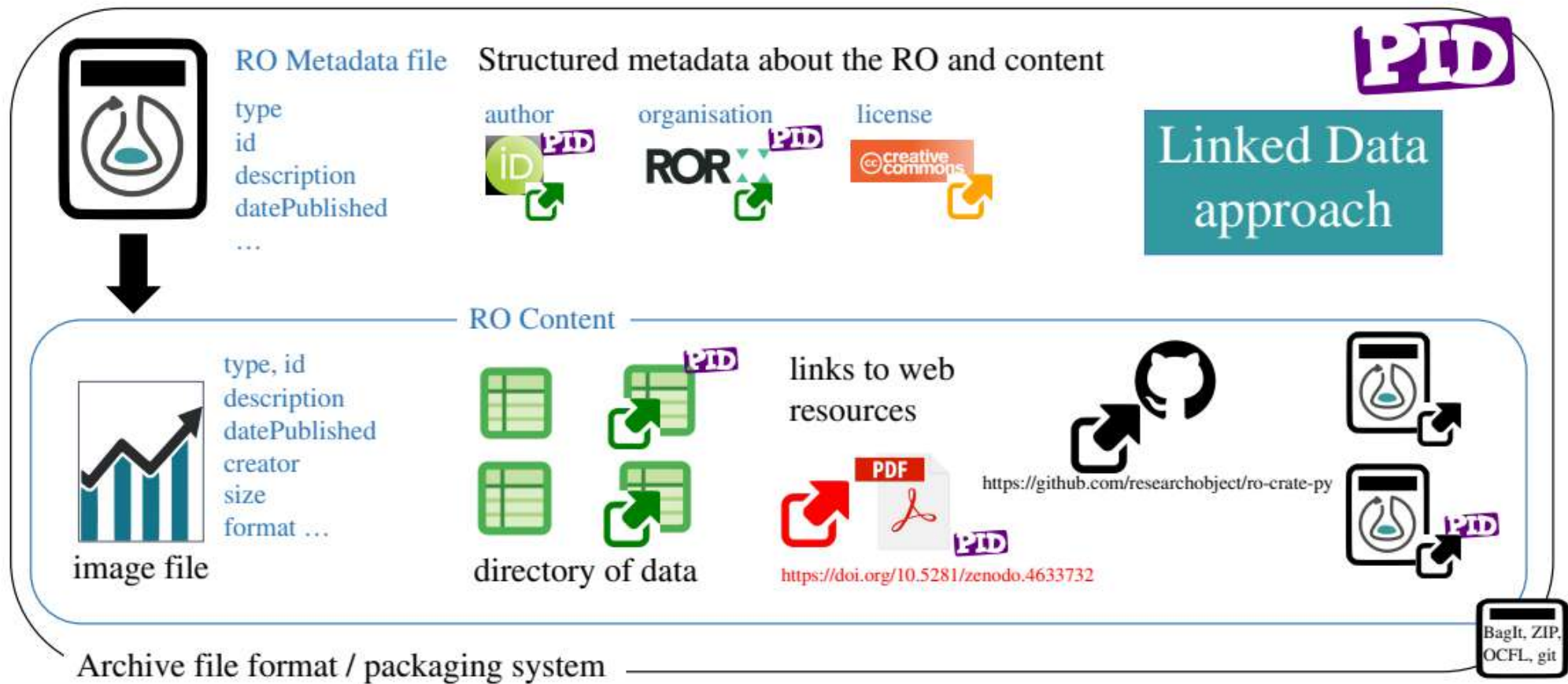


Fig. 1. **Conceptual overview of RO-Crate.** A *Persistent Identifier* (PID) [25] points to a *Research Object* (RO), which may be archived using different packaging approaches like BagIt [26], OCFL [27], git or ZIP. The RO is described within a *RO-Crate Metadata File*, providing identifiers for *authors* using ORCID, *organisations* using Research Organization Registry (ROR) [28] and licences such as Creative Commons using SPDX identifiers. The *RO-Crate content* is further described with additional metadata following a *Linked Data* approach. Data can be embedded files and directories, as well as links to external Web resources, PIDs and nested RO-Crates.

```

{ "@context": "https://w3id.org/ro/crate/1.1/context",
  "@graph": [
    { "@id": "ro-crate-metadata.json",
      "@type": "CreativeWork",
      "conformsTo": {"@id": "https://w3id.org/ro/crate/1.1"},
      "about": {"@id": "."}
    },
    { "@id": "./",
      "@type": "Dataset",
      "name": "A simplified RO-Crate",
      "author": {"@id": "#alice"},
      "license": {"@id": "https://spdx.org/licenses/CC-BY-4.0"},
      "datePublished": "2021-11-02T16:04:43Z",
      "hasPart": [
        {"@id": "survey-responses-2019.csv"},
        {"@id": "https://example.com/pics/5707039334816454031_o.jpg"}
      ]
    },
    { "@id": "survey-responses-2019.csv",
      "@type": "File",
      "about": {"@id": "https://example.com/pics/5707039334816454031_o.jpg"},
      "author": {"@id": "#alice"}
    },
    { "@id": "https://example.com/pics/5707039334816454031_o.jpg",
      "@type": ["File", "ImageObject"],
      "contentLocation": {"@id": "http://sws.geonames.org/8152662/"},
      "author": {"@id": "https://orcid.org/0000-0002-1825-0097"}
    },
    { "@id": "#alice",
      "@type": "Person",
      "name": "Alice"
    },
    { "@id": "https://orcid.org/0000-0002-1825-0097",
      "@type": "Person",
      "name": "Josiah Carberry"
    },
    { "@id": "http://sws.geonames.org/8152662/",
      "@type": "Place",
      "name": "Catalina Park"
    },
    { "@id": "https://spdx.org/licenses/CC-BY-4.0",
      "@type": "CreativeWork",
      "name": "Creative Commons Attribution 4.0"
    }
  ]
}

```



```

{ "@id": "./",
  "@type": "Dataset",
  "name": "A simplified RO-Crate",
  "author": {"@id": "#alice"},
  "license": {"@id": "https://spdx.org/licenses/CC-BY-4.0"},
  "datePublished": "2021-11-02T16:04:43Z",
  "hasPart": [
    {"@id": "survey-responses-2019.csv"},
    {"@id": "https://example.com/pics/5707039334816454031_o.jpg"}
  ]
},
{ "@id": "survey-responses-2019.csv",
  "@type": "File",
  "about": {"@id": "https://example.com/pics/5707039334816454031_o.jpg"},
  "author": {"@id": "#alice"}
},

```

<https://www.researchobject.org/2021-packaging-research-artefacts-with-ro-crate/manuscript.html>

SUMMARY

You should know and be able to explain

- Why we need FAIR principles
- Differences between specific principles
 - Provide your own examples
- Relation between FAIR principles, machine-actionability and open data
- How to apply FAIR principles in practice

Read and watch more about FAIR

Principles explained (by their authors)

- <https://www.nature.com/articles/sdata201618>
- <https://www.go-fair.org/fair-principles/>

Watch (why FAIR matters)

- <https://vimeo.com/143245835>

Related papers

- 'FAIR vs Open' <https://insights.uksg.org/articles/10.1629/uksg.468/>

FAIR cookbook

- <https://faircookbook.elixir-europe.org/content/home.html>

Let's Make Our Data FAIR! Webinar for GO-FAIR

- <https://www.youtube.com/watch?v=dEV2Hnraqal>

FAIR underlies European Open Science Cloud

- <https://eosc-launch.eu/declaration/>

Work with us: TUW Data Services

- Topics:
 - FAIR, PIDs, Data repositories, virtual research environments, Jupyter, Kubernetes, earth observation,...
- Skills
 - 'Coding and setting up things'
- Part time job
 - First contract for 5-6 months
 - Up to 20 hours/week
- Details to be discussed directly with you
- Contact
 - tomasz.miksa@tuwien.ac.at

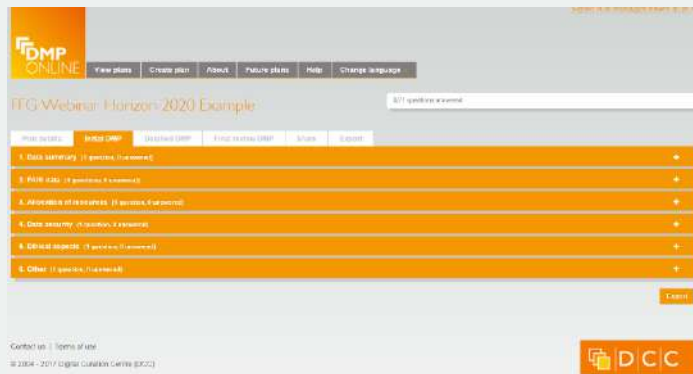
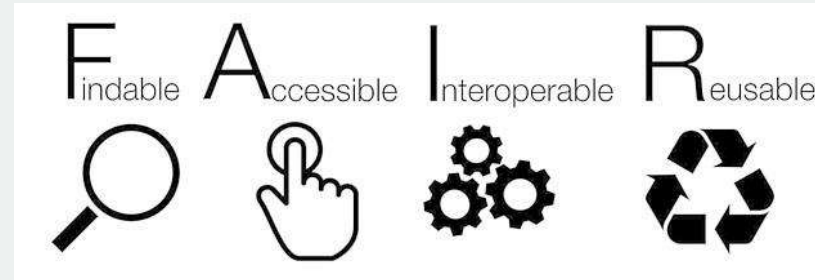
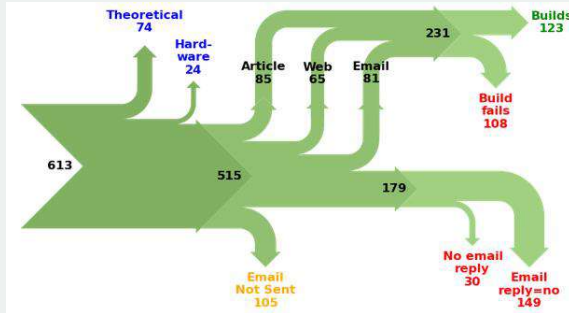
Machine-actionable Data Management Plans

Tomasz Miksa

SBA Research & TU Wien

tmiksa@sba-research.org

Previous lecture on DMPs



Data Management Plans (DMPs) currently

	Data Officer	<i>Who is responsible for the data management and the DMP of the project (name/email address)?</i>
I	Data Characteristics	
I.1	Description of the data	<i>What kinds of data/source code will be generated or reused (type, format, volume)? How will the research data be generated and which methods will be used? How will you structure the data and handle versioning? Who is the target audience?</i>
II	Documentation and Metadata	
II.1	Metadata standards	<i>What metadata standards (if any) will be in use and why? (see Digital Curation Centre)</i>
II.2	Documentation of data	<i>What information is needed for the data to be findable, accessible, interoperable and re-usable (FAIR) in the future? Is the data machine-readable? How are you planning to document this information?</i>
II.3	Data quality control	<i>What quality assurance processes will you adopt? How will the consistency and quality of data collection be controlled and documented? (This may include processes such as repeat samples or measurements, standardised data capture, peer review of data or representation with controlled vocabularies.)</i>
III	Data Availability and Storage	
III.1	Data sharing strategy	<i>How and when will the data be shared and made accessible? What repository will you be using? What persistent identifier will be used?</i>
III.2	Data storage strategy	<i>What data are to be preserved for the long-term, and what data will not be stored? How and where will the data be stored and backed up during the research? How and where will the data be stored after the project ends? For how long will the data be stored? Are there any costs that need to be covered for storage? At what point during or after the project will the data be stored? Are there any technical barriers to making the research data fully or partially accessible?</i>

https://www.fwf.ac.at/fileadmin/files/Dokumente/Open_Access/FWF_DMPTemplate_e.pdf

Data Management Plans (DMPs)

manually created text documents

considered as bureaucracy

created too late

vague

depend on human factor

- scrupulousness
- awareness



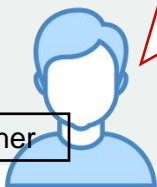
Data Management Plans



How to discover these tools?
Which one do I need to use?
Why do I have to provide the same
information again?

Why haven't they consulted us before?
Who is going to pay for this?
We don't have enough people for that!

Researcher



Stakeholders



Research data lifecycle

Stakeholders involved in research data management

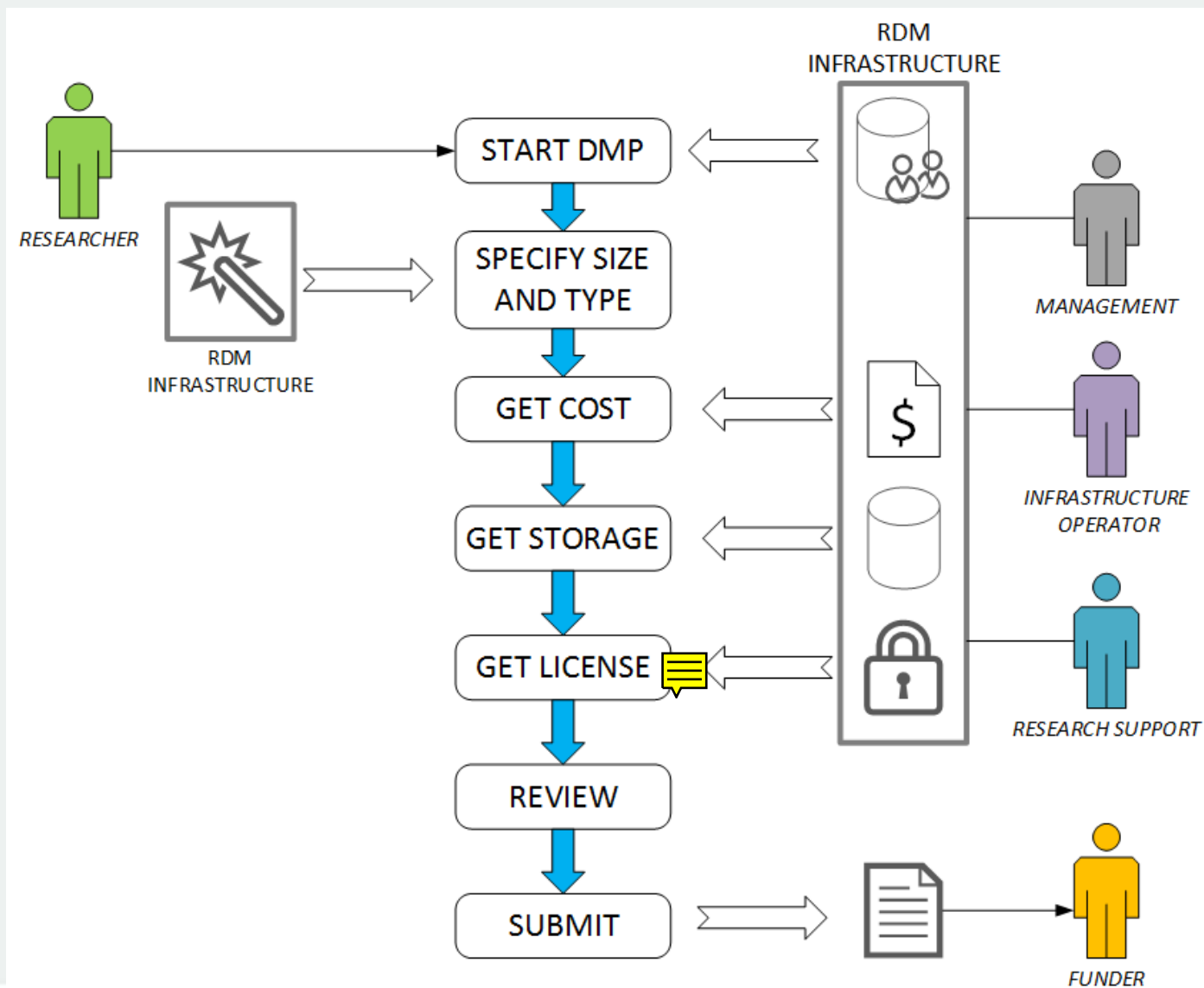
- require information at certain stages
- can provide information if requested at a proper stage

Many problems can be avoided when

- timing is right
- information flow is ensured



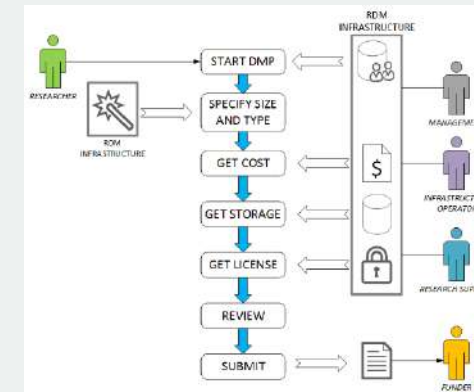
Automated Data Management Workflow



Automated Data Management Workflow

Requires

- common data model
 - to exchange information
- well-defined RDM workflows
 - Who? What? When? How?
- data management infrastructure
 - systems and services implementing workflows



Example

- Current DMPs – model questionnaires

```
<administrative_data>
```

```
  <question>Who will be the Principle Investigator?</question>
```

```
  <answer>The PI will be John Smith from our university.</answer>
```

```
</administrative_data>
```

- Machine-actionable DMPs – model information

```
"dc:creator":[ {
```

```
  "foaf:name":"John Smith",
```

```
  "@id":"orcid.org/0000-1111-2222-3333",
```

```
  "foaf:mbox":"mailto:jsmith@tuwien.ac.at",
```

```
  "madmp:institution":" AT-Vienna-University-of-Technology"
```

```
  } ],
```



Example

- Currently available – not very useful

```
<administrative_data>  
  <question>Who will be the Principle Investigator?</question>  
  <answer>The PI will be John Smith from our university.</answer>  
</administrative_data>
```

- Machine-actionable

Reuse existing standards, e.g. Dublin Core, PREMIS, etc.

```
"dc:creator":[ {  
  "foaf:name":"John Smith",  
  "@id":"orcid.org/0000-1111-2222-3333",  
  "foaf:mbox":"mailto:jsmith@tuwien.ac.at",  
  "madmp:institution":"AT-Vienna-University-of-Technology"  
}],
```

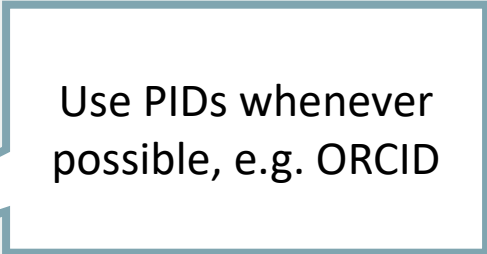

Example

- Currently available – not very useful

```
<administrative_data>  
  <question>Who will be the Principle Investigator?</question>  
  <answer>The PI will be John Smith from our university.</answer>  
</administrative_data>
```

- Machine-actionable DMP

```
"dc:creator":[ {  
  "foaf:name":"John Smith",  
  "@id":"orcid.org/0000-1111-2222-3333",  
  "foaf:mbox":"mailto:jsmith@tuwien.ac.at",  
  "madmp:institution":"AT-Vienna-University-of-Technology"  
}],
```



Use PIDs whenever possible, e.g. ORCID

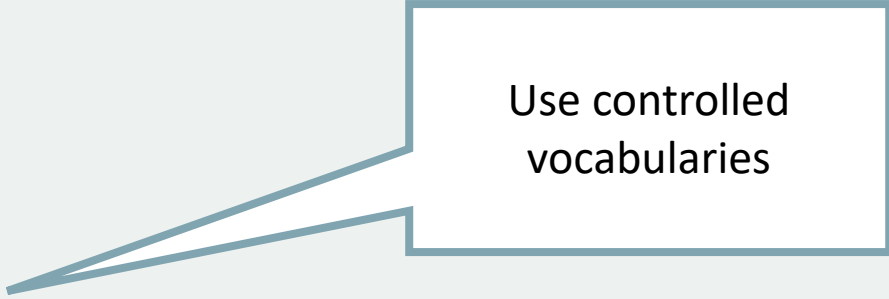
Example

- Currently available – not very useful

```
<administrative_data>  
  <question>Who will be the Principle Investigator?</question>  
  <answer>The PI will be John Smith from our university.</answer>  
</administrative_data>
```

- Machine-actionable DMP

```
"dc:creator":[ {  
  "foaf:name":"John Smith",  
  "@id":"orcid.org/0000-1111-2222-3333",  
  "foaf:mbox":"mailto:jsmith@tuwien.ac.at",  
  "madmp:institution":"AT-Vienna-University-of-Technology"  
}],
```



Use controlled
vocabularies

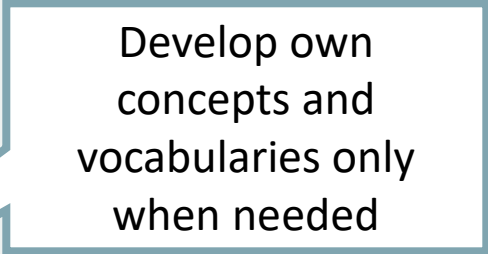
Example

- Currently available – not very useful

```
<administrative_data>  
  <question>Who will be the Principle Investigator?</question>  
  <answer>The PI will be John Smith from our university.</answer>  
</administrative_data>
```

- Machine-actionable DMP

```
"dc:creator":[ {  
  "foaf:name":"John Smith",  
  "@id":"orcid.org/0000-1111-2222-3333",  
  "foaf:mbox":"mailto:jsmith@tuwien.ac.at",  
  "madmp:institution":"AT-Vienna-University-of-Technology"  
}],
```



Develop own
concepts and
vocabularies only
when needed

What is RDA

Research Data Alliance

- community-driven organization
- 6,000 members from 130 countries
- different stakeholders

Plenary meetings

Interest Groups (IGs)

- Active DMPs

Working Groups (WGs)

- DMP Common Standards



<https://rd-alliance.org/>

DMP Common Standards WG

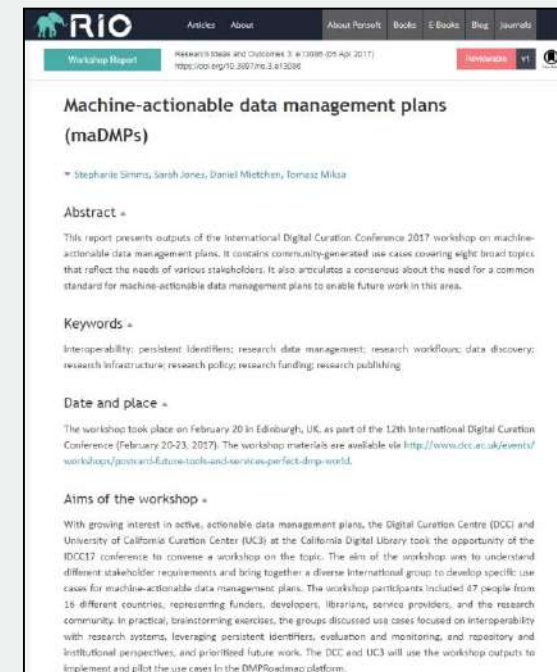
Launched in October 2017

Result of a consultation made by Active DMPs IG

Focus on machine-actionable DMPs

280 members from all continents

DMP tool owners are part of it



<https://doi.org/10.3897/rio.3.e13086>

DMP Common Standards WG

Taxonomy:



Posts



Create Wiki index



Events



Repository



Outputs



Case Statements



Plenaries



Members

[create new content](#) ▾Group Status:  WGs Maintaining deliverables (maintenance group)[Join Group](#)**Status:** Recognised & Endorsed**Chair (s):** Paul Walk, Peter Neish, Tomasz Miksa**Secretariat Liaison:** enquiries[at]rd-alliance.org**TAB Liaison:** Isabelle Perseil

The need for establishing this working group was articulated during the 9th plenary meeting in Barcelona during the Active DMPs IG session. The discussion was framed by a white paper by Simms et al. on machine-actionable data management plans (DMPs). The white paper is based on outputs from the IDCC workshop held in Edinburgh in 2017 that gathered almost 50 participants from Africa, America, Australia, and Europe. It describes eight community use cases which articulate consensus about the need for a common standard for machine-actionable DMPs (where machine actionable is defined as "information that is structured in a consistent way so that machines, or computers, can be programmed against the structure")

The specific focus of this working group is on developing common information model and specifying access mechanisms that make DMPs machine-actionable. The outputs of this working group will help in making systems interoperable and will allow for automatic exchange, integration, and validation of information provided in DMPs, for example, by checking whether a provided PID links to an existing dataset, if hashes of files match to their provenance traces, or whether a license was specified. The common information models are NOT intended to be prescriptive templates or questionnaires, but to provide re-usable ways of representing machine-actionable information on themes covered by DMPs.

<https://www.rd-alliance.org/groups/dmp-common-standards-wg>

SCOPING MADMPS

Scoping maDMPs by DMP Common Standards WG

1st consultation

2nd consultation

Proof of concept tools

BPMN processes

Model development

1st consultation – user stories

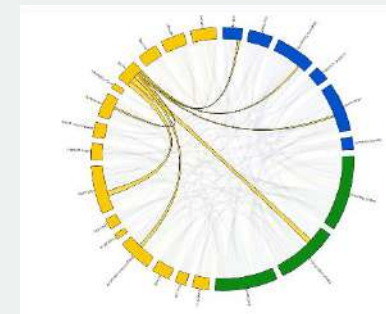
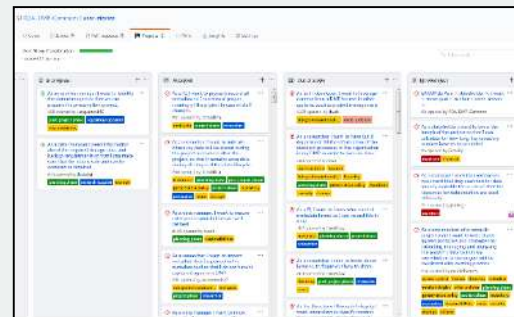
Goals

- identify stakeholders at each lifecycle stage
 - define which information they **provide**
 - define which information they **expect**



As a <stakeholder>, I want <goal> so that <reason >.

*As a **researcher**, I want to **inform repository operator** on the amount of data in the planning phase, so that they provide **information on costs**.*



2nd consultation – existing models

2nd consultation goes deep

- how do we model specific requirements
 - which specific fields are needed?
 - which models exist?



(Meta-) Data	Requirements
<h3>Overview</h3> <p>This document is part of a consultation described here: link</p> <p>From the previous consultation with user stories we have derived following high level requirements:</p> <ul style="list-style-type: none">• Format<ul style="list-style-type: none">◦ Format [80, 12, 99, 62, 67, 54, 80]• Volume<ul style="list-style-type: none">◦ Data size estimate [5, 77, 80, 100]<ul style="list-style-type: none">▪ For specific type of data [62]◦ Data size real [34]• Provenance [54]• Metadata<ul style="list-style-type: none">◦ taxonomy/classification [14, 11]◦ Links to metadata of the real data [89, 39]◦ Link publications to data [55]◦ Authorship [88]◦ Multilingual metadata [65]◦ Include raw metadata directly in the model [91, 65]• Reuse<ul style="list-style-type: none">◦ Links to (meta-)data location [89, 90, 96, 39, 00]• Repository [42]<ul style="list-style-type: none">◦ Persistent identifier for data [92]◦ Link publications to data [55, 88]◦ Link to License/Contract allowing data usage/storing [56] <p>Please help us:</p> <ul style="list-style-type: none">• Break down existing requirements into more specific requirements,• Add missing requirements,• Provide examples of existing models, vocabularies, etc. that can be used to model these <p>Please provide your suggestions below.</p>	<p>Quality - <code>dqv:hasQualityAnnotation</code> (statement related to quality of the Dataset, including rating, quality certificate, feedback that can be associated to the Dataset). <code>Stat:dimension</code>, <code>stat:measure</code></p> <p>Data Dimensions and units of measurement (<code>stat:dimension</code>, <code>stat:measure</code>)</p> <h3>Models</h3> <p>Format <code>dct:format</code></p> <p>Volume <code>dct:generalPeriodicity</code></p> <p>Provenance <code>dct:creator</code>, <code>dcat:contactPoint</code>, <code>prov:generated</code>, <code>prov:qualifiedAttribution</code></p> <p>Metadata Taxonomy/classification - <code>dct:subject</code>, <code>dcat:theme</code> Link publication to data: <code>dct:relations</code> (link to Publications catalogue), <code>adms:identifier</code> (link to related publication-identifiers such as DOI, ISSN, ISBN) Authorship - <code>dct:publisher</code>, <code>prov:agent</code>, <code>foaf:name</code> Conformity to data model - <code>dct:conformsTo</code> Multilingual metadata - <code>dct:language</code> Include raw data in the data model - <code>adms:sample</code> (refers to a sample of data)</p> <p>Reuse Links to metadata location - <code>dct:source</code>, <code>foaf:homepage</code> (documentation)</p> <p>Repository Persistent identifier for data - <code>dct:identifier</code> Link publications to data - <code>dcat:distribution</code> License/contract - <code>dct:accessRights</code>, <code>dct:license</code></p> <h3>Other comments</h3> <p>https://joinup.ec.europa.eu/release/statdcat-ep-v100 https://joinup.ec.europa.eu/release/statdcat-ep-v11</p>

Proof of concept tools

Requirements

- Provide minimum input
- Import as much as possible from existing systems to help in creating maDMPs

Tools available as Docker on GitHub

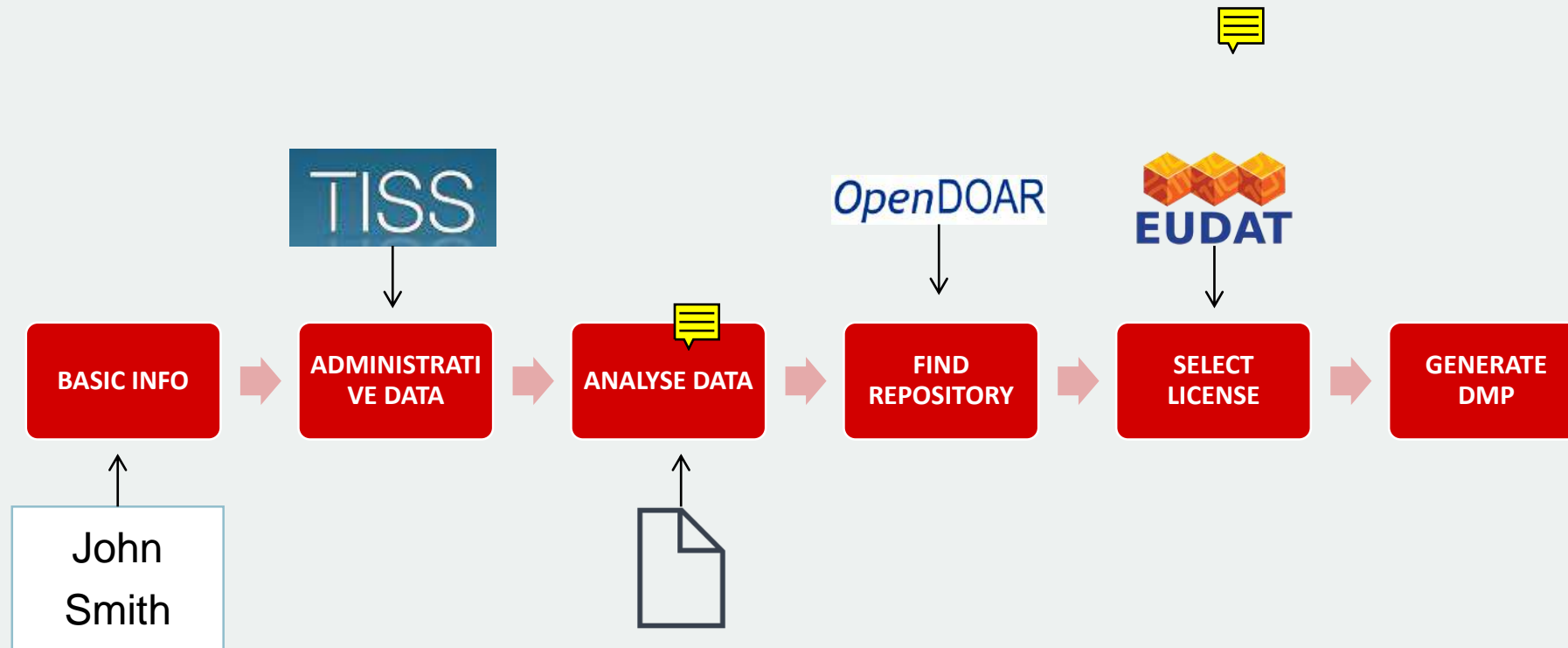
- <https://github.com/TomMiksa/DMPGenerator>
- https://github.com/TomMiksa/digital_preservation_ex_1_2
- <https://github.com/TomMiksa/tu-dpue-lab2-ss18>
- https://github.com/TomMiksa/DigitalPreservation_2
- <https://github.com/TomMiksa/digitalpreservation-dmp-generator>
- <https://github.com/TomMiksa/DMPPlanner>

Example of a landing page for maDMPs

- <https://oblassers.github.io/fair-data-science/>
- <https://github.com/oblassers/fair-data-science>

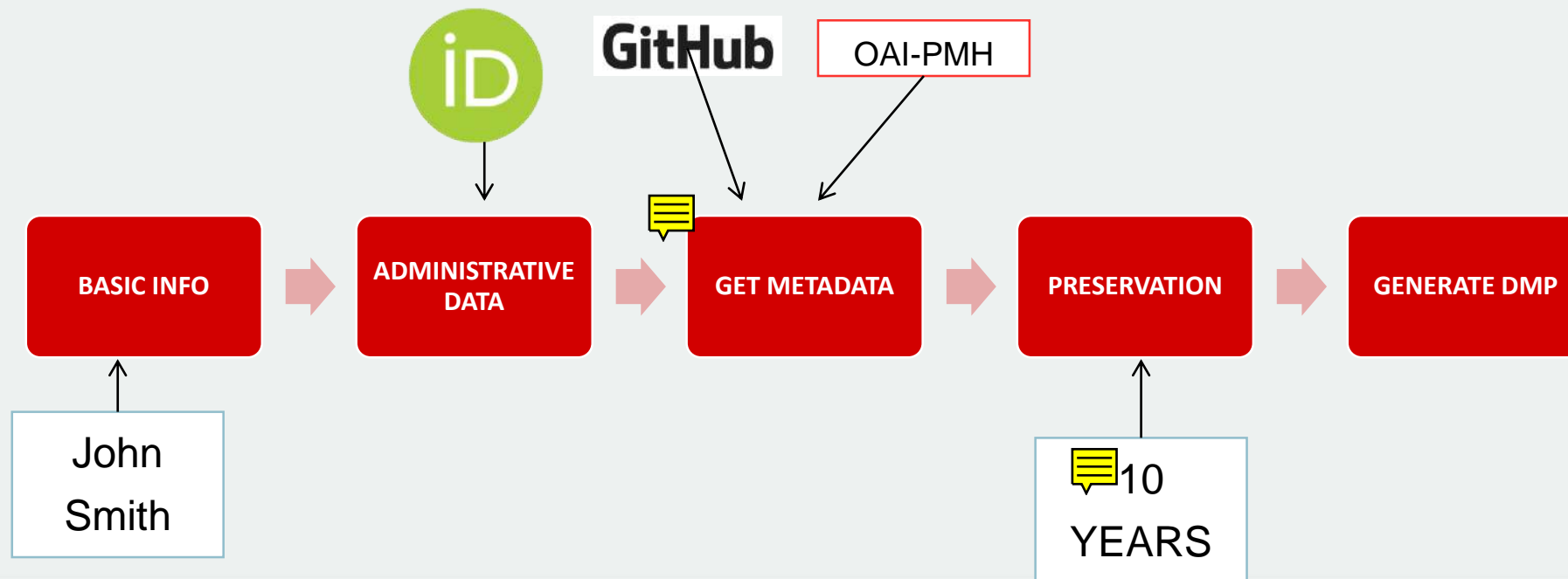
Planning phase

- Goal: get **estimations** and **recommendations** (which are feasible to implement later)



Project and Post-project phases

- Goal: **update** DMP with **real** information by **re-using** (linking) information provided elsewhere





Name

Please provide your full name.

full_name Tomasz Miksa

orcid 0000-0002-4929-7875

current_employment_name SBA Research



Resources

Add as many Github repositories or OAI-PMH compliant DOIs as you like.*

Zenodo Ten Simple Rules For Machine-Actionable Data Management Plans (Preprint)

documentation

Remove

Github TomMiksa/DMPlanner

software

Remove



Preservation Time

Choose how many years the data for each group should be kept.

Software 10 years

Documentation 20 years

TUW DMP

A Data Management Plan created using DMPlanner.

Creator

Name: Tomasz Miksa

ORCID: [0000-0002-4929-7875](https://orcid.org/0000-0002-4929-7875)

Current Work: SBA Research

How will you manage copyright and Intellectual Property Rights (IPR) issues?

The software which was created in the course of the project has the license restrictions "MIT License".

Which data are of long-term value and should be retained, shared, and/or preserved?

In this project especially the documentation, as well as the software has a long-term value and should at least be as long preserved as the targeted preservation time specifies. The targeted preservation time for the documentation is 20 years. The targeted preservation time for the software is 10 years.

What is the long-term preservation plan for the dataset?

One of the main strategies of the long-term preservation plan is the use of public accessible repositories to save the components of the project. The documentation resource "Ten Simple Rules For Machine-Actionable Data Management Plans (Preprint)" is hosted on Zenodo. The software resource "DMPlanner" is hosted on Github.

How will you share the data?

The data will be primarily shared through the public repositories listed above. This way the data is openly accessible and findable, as well as searchable. The data is available at the repositories as of this moment.

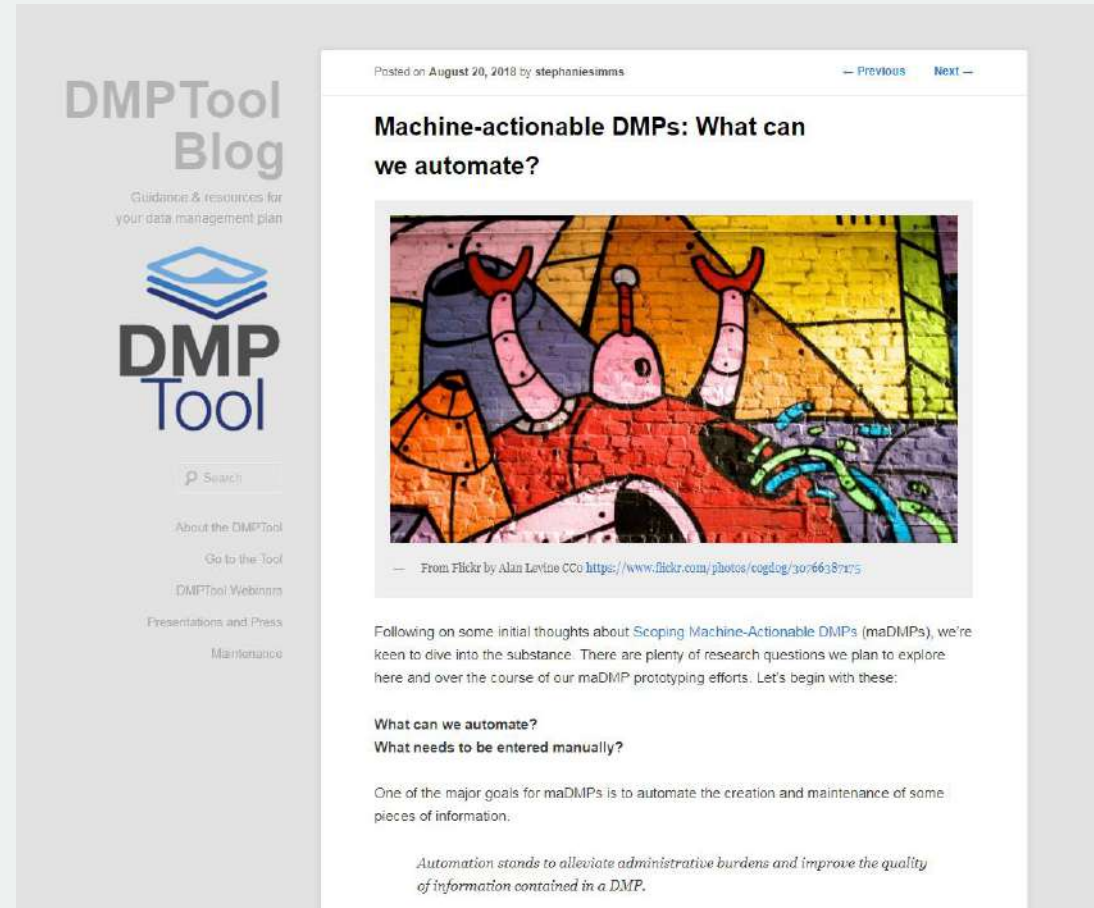
Are any restrictions on data sharing required?

The restrictions on data sharing are composed of the used licenses together with the long-term preservation plan. With this in mind the following restrictions for the resources of the project apply. The documentation resource "Ten Simple Rules For Machine-Actionable Data Management Plans (Preprint)" will be hosted on Zenodo for at least 20 years. The software resource "DMPlanner" will be hosted on Github for at least 10 years.

Who will be responsible for data management?

The creator of this data management plan is Tomasz Miksa. Therefore Tomasz Miksa is also the reference person for possible reviews and revisions regarding this data management plan in the future. Unless amended Tomasz Miksa is additionally responsible for the adherence to the plan.


Machine-actionable DMPs: What can we automate?



The image shows a screenshot of a blog post from the DMPTool Blog. The page layout includes a sidebar on the left with the DMPTool logo and navigation links, and a main content area on the right. The main content area features a post titled "Machine-actionable DMPs: What can we automate?" by stephaniesimms, dated August 20, 2018. The post includes a colorful abstract illustration of a robot-like figure and a list of research questions to explore.

Posted on August 20, 2018 by stephaniesimms — Previous Next —

Machine-actionable DMPs: What can we automate?



— From Flickr by Alan Levine CCo <https://www.flickr.com/photos/cogdog/30766387175>

Following on some initial thoughts about [Scoping Machine-Actionable DMPs \(maDMPs\)](#), we're keen to dive into the substance. There are plenty of research questions we plan to explore here and over the course of our maDMP prototyping efforts. Let's begin with these:

What can we automate?
What needs to be entered manually?

One of the major goals for maDMPs is to automate the creation and maintenance of some pieces of information.

Automation stands to alleviate administrative burdens and improve the quality of information contained in a DMP.

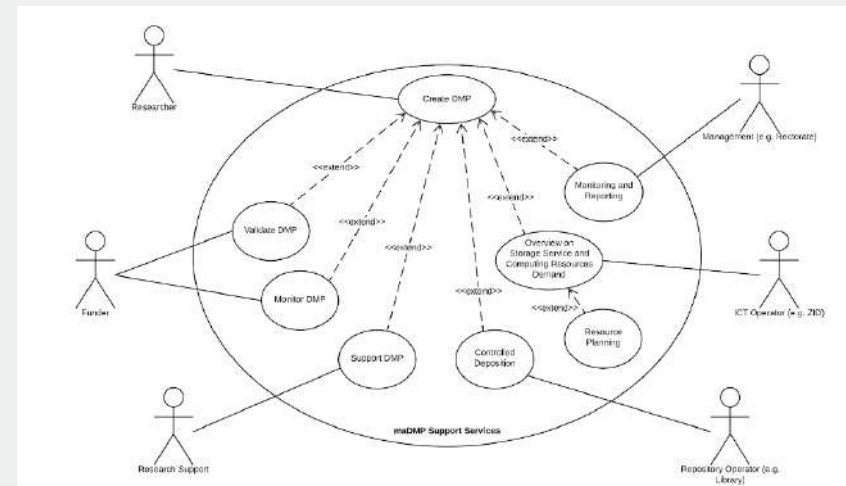
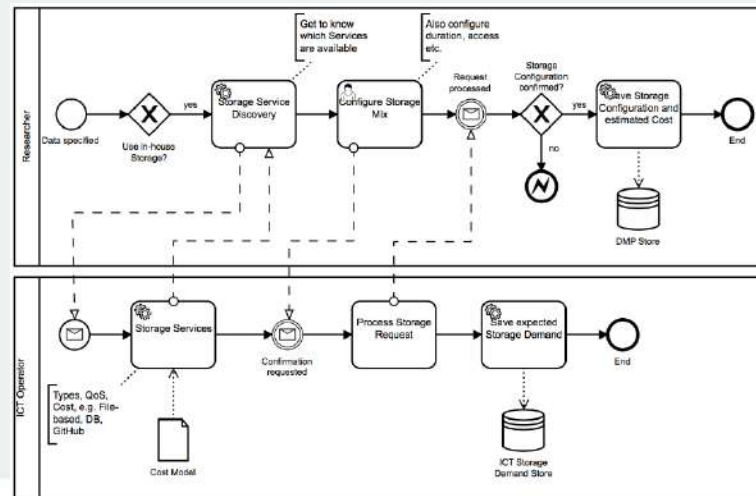
<https://blog.dmptool.org/2018/08/20/machine-actionable-dmps-what-can-we-automate/>

<https://blog.dmptool.org/2019/10/14/whats-new-with-our-machine-actionable-dmp-work/>

Processes

Processes help identify

- **tasks** performed by stakeholders
 - e.g. ICT operator provide costs of storage
- **systems** needed to be put in place
 - e.g. maDMP repository or costing service
- **concepts** to be developed or agreed
 - e.g. cost model for storage



Processes

Useful in deploying maDMPs

Allow us to narrow down focus

- common model does not contain business logic
 - e.g. cost estimation is done by a service that provides a value
- common model is an information carrier
 - tools, services, processes make maDMPs *machine-actionable*

BPMN Processes for machine-actionable DMPs

Oliver Ottaviani & Tommaso Milani

Contents	2
Start DMP	2
Specify Service Type	2
Get Cost and Ratings	4
Request Configuration and Cost Estimation	4
Request Provisioning	5
Get Results	5
Get Metadata/Statistics	7
Get Availability	8
Cancel Order	8
Get Help	10

<http://doi.org/10.5281/zenodo.2607556>

Scoping maDMPs - summary

1st consultation (user stories) went broad

- to define scope of maDMPs

2nd consultation went deep

- to identify models for specific requirements

Proof of concept tools

- to demonstrate how model can be used to automate tasks

BPMN processes

- to identify systems and stakeholders involved

Model development

Official RDA Recommendation on maDMPs



RDA DMP Common Standard for Machine-actionable Data Management Plans

The Challenge:

Data Management Plans are free-form text documents describing the data that is used and produced during the course of research activities. They specify where the data will be archived, which licenses and constraints apply, and to whom credit should be given, etc. The workload and bureaucracy often associated with traditional DMPs can be reduced when they become machine-actionable.



Produced by: **DMP Common Standards WG**
<https://www.rd-alliance.org/groups/dmp-common-standards-wg>

RDA DMP Common Standard for Machine-actionable Data Management Plans

Recommendations of the RDA DMP Common Standards WG
Tomasz Mikso, Paul Walk, Peter Neish

Purpose

This application profile is meant for exchange of machine-actionable DMPs between systems. It is independent of any internal data organisation used by these systems. The application profile does not prescribe how information must be presented to the end user and does not enforce any specific logic on how this information must be collected or used. The application profile is an information carrier and the full machine-actionability can only be achieved when systems using the application profile implement appropriate logic.

This application profile is intended to cover a wide range of use cases and does not set any business (e.g. funder specific) requirements. It represents information over the whole DMP lifecycle, that is, it can express planned actions, as well as actions already performed.

The application profile is NOT intended to be a prescriptive template or a questionnaire, but to provide a re-usable way of representing machine-actionable information on themes covered by DMPs.

Overview

Figure 1 presents concepts used within the application profile. Each concept is further broken down into specific fields (not depicted). The full application profile specification can be found [online](#). Below we outline main concepts used within the application profile that are depicted in Figure 1.

DMP - Provides high level information about the DMP, e.g. its title, modification date, etc. It is the root of this application profile.

Project - Describes the project associated with the DMP, if applicable. It can be used to describe any type of project: that is, not only funded projects, but also internal projects, PhD theses, etc.

Funding - For specifying details on funded projects, e.g. NSF or EC funded projects.

Contact - Specifies the party which can provide information on the DMP.

Contributor - For listing all parties involved in the process of data management described by DMPs.

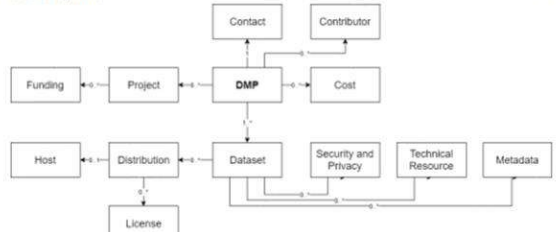


Figure 1: Overview of the application profile for the machine-actionable DMPs.

1

Miksa, T., Walk, P., & Neish, P. (2020). RDA DMP Common Standard for Machine-actionable Data Management Plans. <https://doi.org/10.15497/rda00039>

Pending adoptions (selected)



Webinar available

**HIGHLIGHTING SOLUTIONS PROPOSED BY
RDA ACTIVE DMPS, EXPOSING DMPS AND DMP
COMMON STANDARDS WORKING GROUPS**

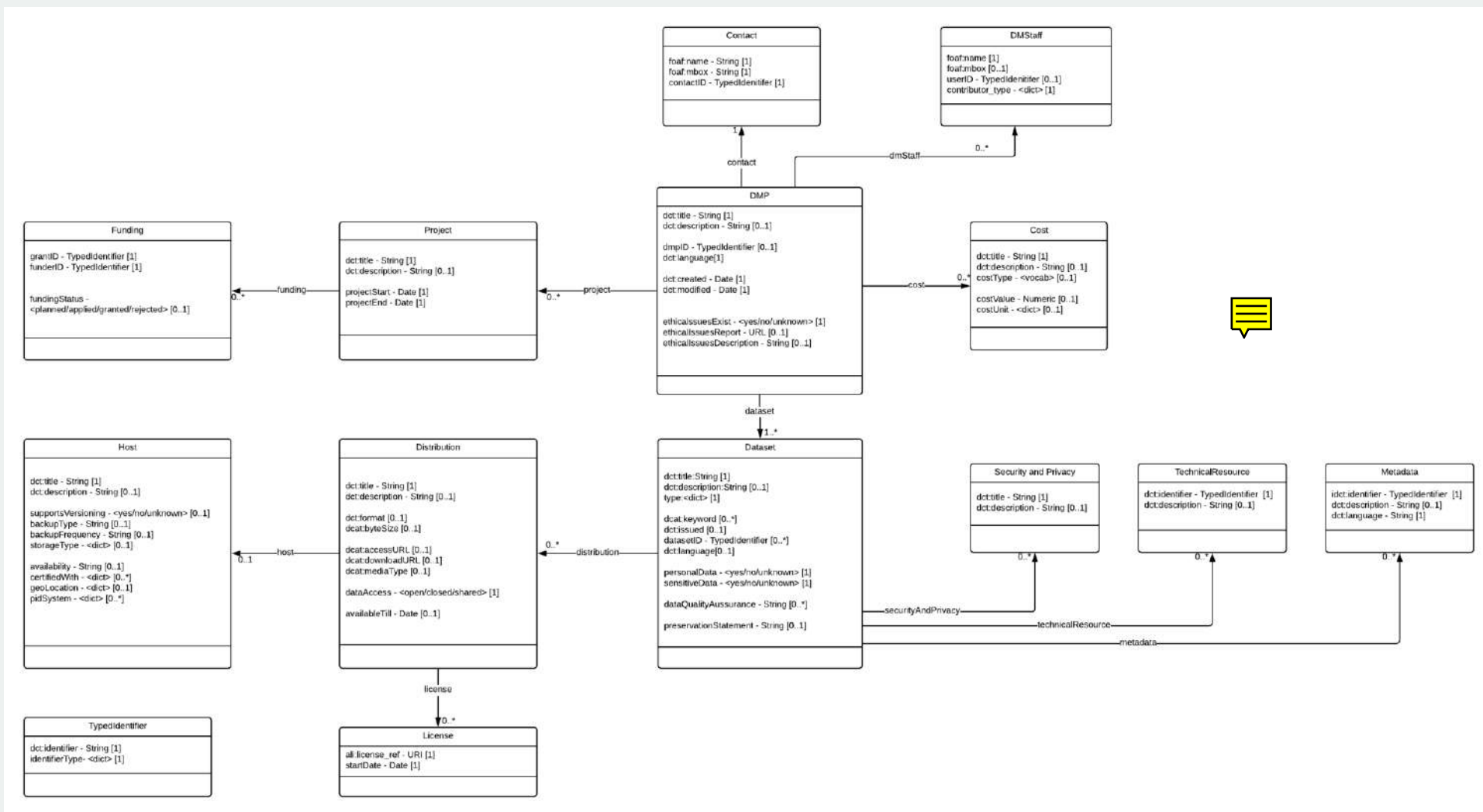
	Daniel Bangert Göttingen State and University Library RDA Secretariat		Kathryn Unsworth Commonwealth Scientific and Industrial Research Organisation RDA Exposing DMPs WG
	Peter Neish University of Melbourne RDA DMP Common Standards WG		Tomasz Miksa SBA Research RDA DMP Common Standards WG
	Sam Rust Digital Curation Centre	<p>WEBINAR SLIDES & RECORDINGS NOW AVAILABLE</p> 	

<https://www.rd-alliance.org/rda-working-groups-solutions-dmp-recording-and-slides-webinar-now-available-0>

Part 2

COMMON STANDARD FOR MADMPS

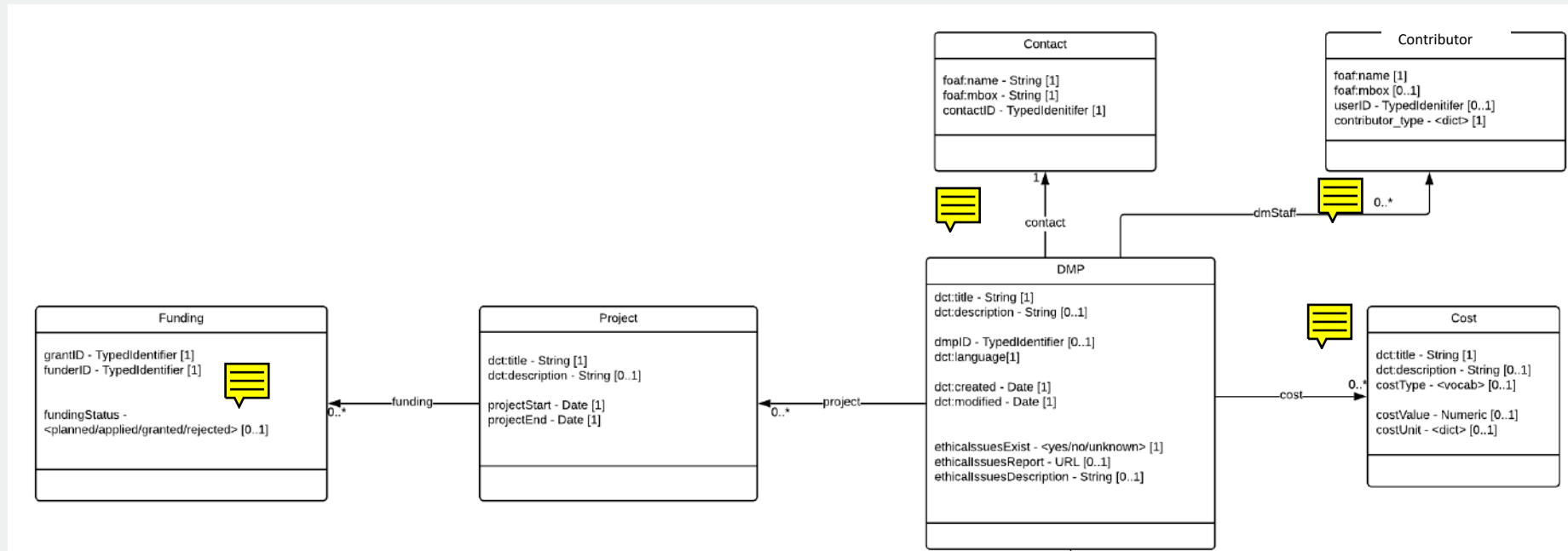
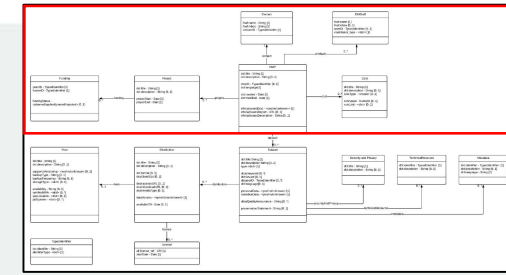
Common standard for maDMPs



Note: Diagram depicts draft version. Details evolved! Use official release.

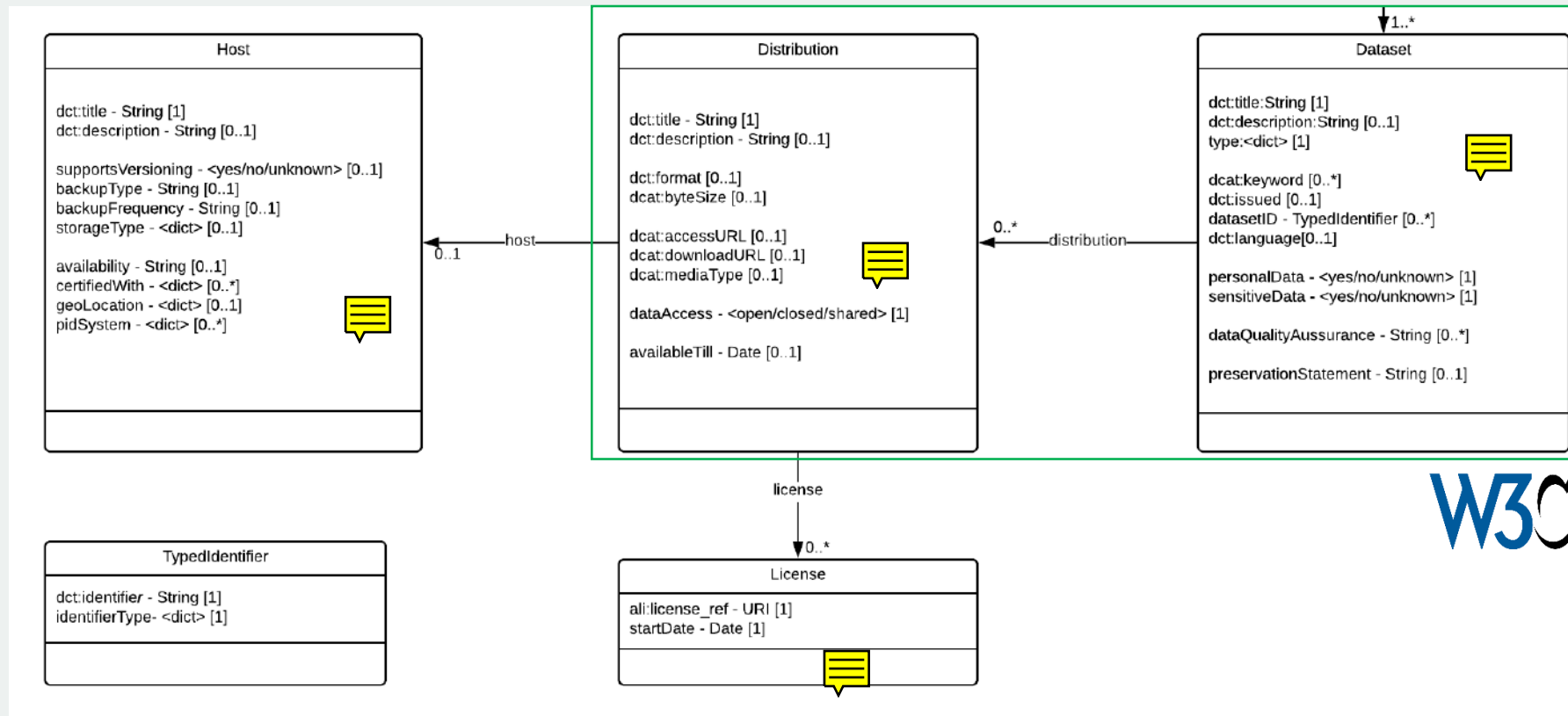
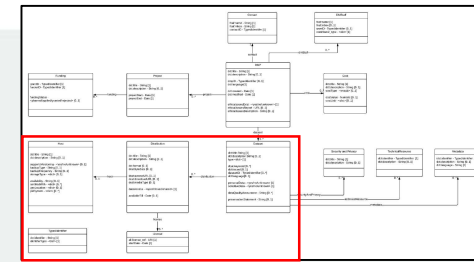
<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard>

Common standard for maDMPs



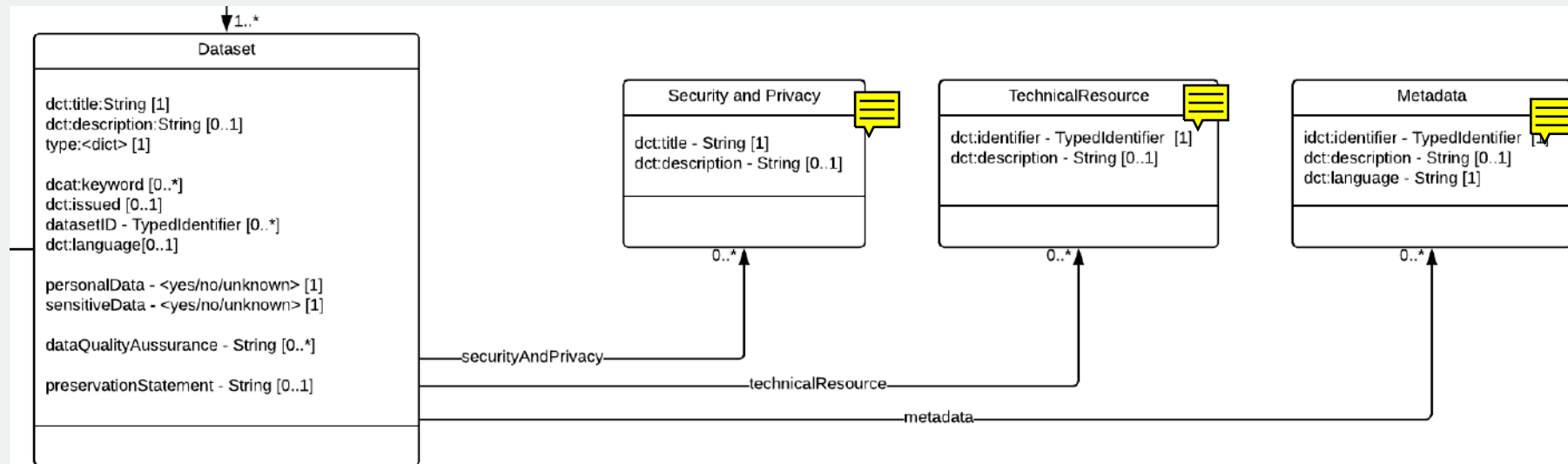
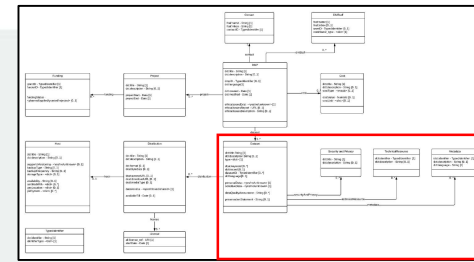
Note: Diagram depicts draft version. Details evolved! Use official release.

Common standard for maDMPs



Note: Diagram depicts draft version. Details evolved! Use official release.

Common standard for maDMPs



Note: Diagram depicts draft version. Details evolved! Use official release.

Standard - documentation

Properties in 'contact'

Name	Description	Data Type	Cardinality	Example Value
contact_id	Identifier for a contact person	String	Exactly One	http://orcid.org/0000-0000-0000-0000
mail	E-mail address	String	Exactly One	cc@example.com
name	Name of the contact person	String	Exactly One	Charlie Chaplin

Properties in 'cost'

Name	Description	Data Type	Cardinality	Example Value
currency_code	Allowed values defined by ISO 4217.	Term from Controlled Vocabulary	Zero or One	EUR
description	Description	String	Zero or One	Costs for maintaining....
title	Title	String	Exactly One	Storage and backup
type	Type	Term from Controlled Vocabulary	Zero or One	
value	Value	Number	Zero or One	1000

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/docs/index.md>

Standard – FAQ

The screenshot shows a GitHub repository page for 'RDA-DMP-Common / RDA-DMP-Common-Standard'. The repository has 3 unwatched items, 0 stars, and 5 forks. The current view is for the file 'FAQ.md' in the 'docs' directory, on the 'master' branch. A commit by TomMiksa is shown, updating the file 21 hours ago. The file is 8.34 KB and contains 85 lines of code. The content of the file is as follows:

Frequently Asked Questions

Index:

- [When to use the model?](#)
- [Do I need to populate all fields?](#)
- [What is the granularity of a Dataset?](#)
- [What is a difference between Dataset and a Distribution?](#)
- [How versioning works?](#)
- [How to express something is planned?](#)
- [How to indicate actions that were performed?](#)
- [How to model embargoes?](#)
- [Why Metadata is referenced from a Dataset?](#)
- [Are there any other serialisations planned different than JSON?](#)
- [Is there a JSON Schema?](#)
- [Is there a model validator?](#)

When to use the model?

The model is meant for exchange of machine-actionable DMPs between systems. The model is independent of any internal

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/docs/FAQ.md>

Standard – useful links

The screenshot shows a GitHub repository page for 'RDA-DMP-Common / RDA-DMP-Common-Standard'. The file 'docs / links.md' is selected, showing a commit by TomMiksa from 2 days ago. The file content is as follows:

69 lines (45 sloc) | 3.84 KB

Links

We have collected here links to all important resources created by the [RDA DMP Common Standards WG](#) (official website).

1st Consultation - scoping the maDMPs

Collection of user stories to identify scope of maDMPs.

- [Description of the consultation](#)
- [User stories organised on a project board](#)
- [Interactive visualisation of user stories](#)
- [Report from Vienna workshop for collecting user stories](#)
- [iPres conference paper summarising the consultation](#)

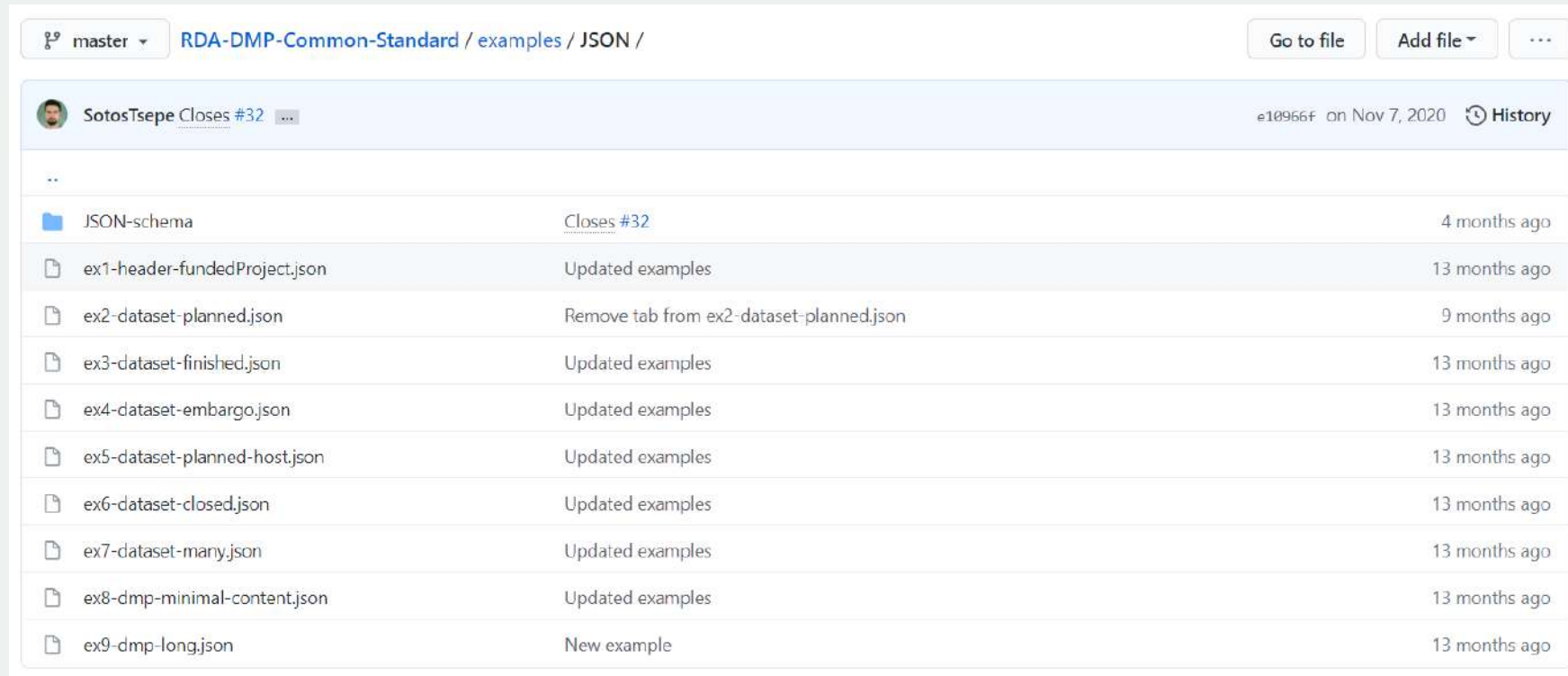
2nd Consultation - existing models

Collection of models that are relevant in view of requirements derived from the user stories

- [Description of the 2nd consultation \(includes further links\)](#)

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/docs/links.md>

Standard – JSON examples



The screenshot shows a GitHub repository interface for the path `RDA-DMP-Common-Standard / examples / JSON /`. The repository is on the `master` branch. A commit by `SotosTsepe` is shown, closing issue #32 on Nov 7, 2020. Below the commit, a list of files and folders is displayed:

File/Folder	Commit Message	Time
..		
JSON-schema	Closes #32	4 months ago
ex1-header-fundedProject.json	Updated examples	13 months ago
ex2-dataset-planned.json	Remove tab from ex2-dataset-planned.json	9 months ago
ex3-dataset-finished.json	Updated examples	13 months ago
ex4-dataset-embargo.json	Updated examples	13 months ago
ex5-dataset-planned-host.json	Updated examples	13 months ago
ex6-dataset-closed.json	Updated examples	13 months ago
ex7-dataset-many.json	Updated examples	13 months ago
ex8-dmp-minimal-content.json	Updated examples	13 months ago
ex9-dmp-long.json	New example	13 months ago

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/tree/master/examples/JSON>

DMP and Project – JSON example

```
40 lines (34 sloc) 825 Bytes Raw Blame History
1 {
2   "DMP": {
3     "title": "Funded DMP",
4     "description": "Example of a DMP header for a funded project.",
5
6     "created": "2019-02-22T13:20:15.5",
7     "modified": "2019-02-22T15:10:56.9",
8     "contact": {
9       "name": "First Last",
10      "inbox": "test@test",
11      "contactID": {
12        "identifier": "https://orcid.org/0000-0002-4929-7875",
13        "identifierType": "HTTP-ORCID"
14      }
15    },
16    "ethicalIssuesExist": "false",
17
18    "project": {
19      "title": "Making maDMPs awesome",
20      "projectStart": "2017-01-01",
21      "projectEnd": "2020-12-31",
22
23      "funding": {
24        "funderID": {
25          "identifier": "501100002428",
26          "identifierType": "FUNDREF"
27        },
28        "grantID": {
29          "identifier": "1234567-AT",
30          "identifierType": "custom"
31        },
32        "fundingStatus": "granted"
33      }
34    },
35
36    "dataset": {}
37  }
38 }
39 }
```

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/examples/JSON/ex1-header-fundedProject.json>

DMP and Project – JSON example

40 lines (34 sloc) | 825 Bytes

```
1  {
2      "DMP": {
3          "title": "Funded DMP",
4          "description": "Example of a DMP header for a funded project.",
5
6          "created": "2019-02-22T13:20:15.5",
7          "modified": "2019-02-22T15:10:56.9",
8          "contact": {
9              "name": "First Last",
10             "mbox": "test@test",
11             "contactID": {
12                 "identifier": "https://orcid.org/0000-0002-4929-7875",
13                 "identifierType": "HTTP-ORCID"
14             }
15         },
16         "ethicalIssuesExist": "false",
```

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/examples/JSON/ex1-header-fundedProject.json>

DMP and Project – JSON example

```
18         "project": {
19             "title": "Making maDMPs awesome",
20             "projectStart": "2017-01-01",
21             "projectEnd": "2020-12-31",
22
23             "funding": {
24                 "funderID": {
25                     "identifier": "501100002428",
26                     "identifierType": "FUNDREF"
27                 },
28                 "grantID": {
29                     "identifier": "1234567-AT",
30                     "identifierType": "custom"
31                 },
32                 "fundingStatus": "granted"
33             }
34         },
35
36         "dataset" : {}
37
38     }
39 }
```

<https://github.com/RDA-DMP-Common/RDA-DMP-Common-Standard/blob/master/examples/JSON/ex1-header-fundedProject.json>

Standard assumptions – relaxed constraints

Model must be applicable in different settings

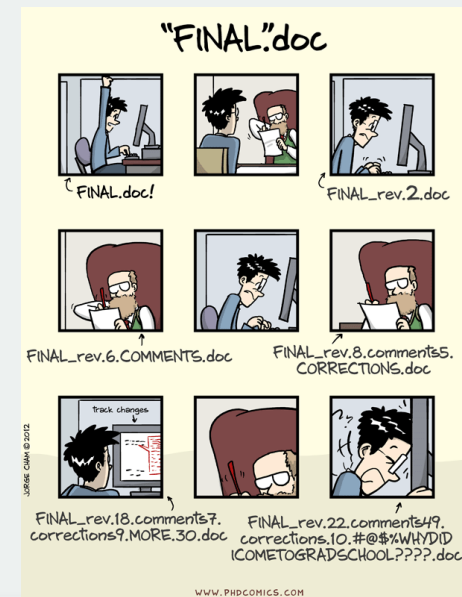
- relaxed constraints within the model
 - e.g. DMP **can** relate to a project [0..*]
- constraints introduced at the ‘business level’
 - tool implementing the model
 - e.g. DMP **must** relate to a project [1..*]
- DMP instances are still compatible



Standard assumptions - interoperability

Model will be pre-dominantly used to exchange information between systems

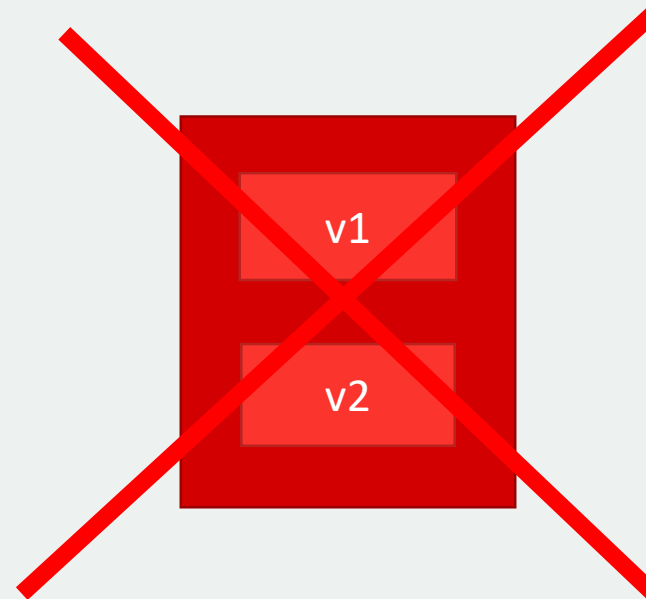
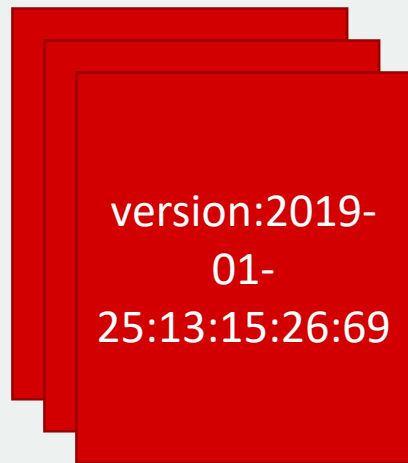
- internal representation of information in a DMP tool may differ (physical model)
 - e.g. database may have a different schema
- No 'meta-fields' about DMP
 - e.g. no DMP state field 'final'



Standard assumptions - versioning

DMP versioning done by systems using the model

- model provides fields allowing to identify DMP version
- model does not track connections between versions



Standard assumptions – evolving information

Model expresses ‘certainty’ of provided information

- to support different phases of DMPs

Example

- Source code will be issued on 2019-06-30 (planned) in ‘some-repo’.
- There will be an embargo period till 2019-12-31.
- Later on the source code will be available on a CC-BY license.

```
"DMP": {  
  "modified": "2019-02-22T13:20:15.5"  
  "dataset": {  
    "title": "Source Code",  
    "issued": "2019-06-30",   
    "distribution": {  
      "accessURL": "http://some-repo...",  
      "license": {  
        "license_ref": "https://creativecommons.org/licenses/by/4.0/",  
        "startDate": "2019-12-31"  
      }  
    }  
  }  
}
```

Standard – reused standards

id	label	uri
ali	Access License and Indicators	http://www.niso.org/schemas/ali/1.0/
dces	Dublin Core Element Set	http://purl.org/dc/elements/1.1/
dct	DCMI Metadata Terms	http://purl.org/dc/terms/
foaf	Friend of a Friend (FOAF)	http://xmlns.com/foaf/0.1/
dcat	DCAT	https://www.w3.org/TR/vocab-dcat/
datacite	Data Cite	https://schema.datacite.org
cerif	Cerif	https://www.eurocris.org/ontologies/cerif/1.3/index.html#currencyCode
coar	COAR	http://vocabularies.coar-repositories.org/pubby/resource_type.html
iso6391	ISO 6391-1	Two letter country code
iso4217	ISO 4217	Currency code

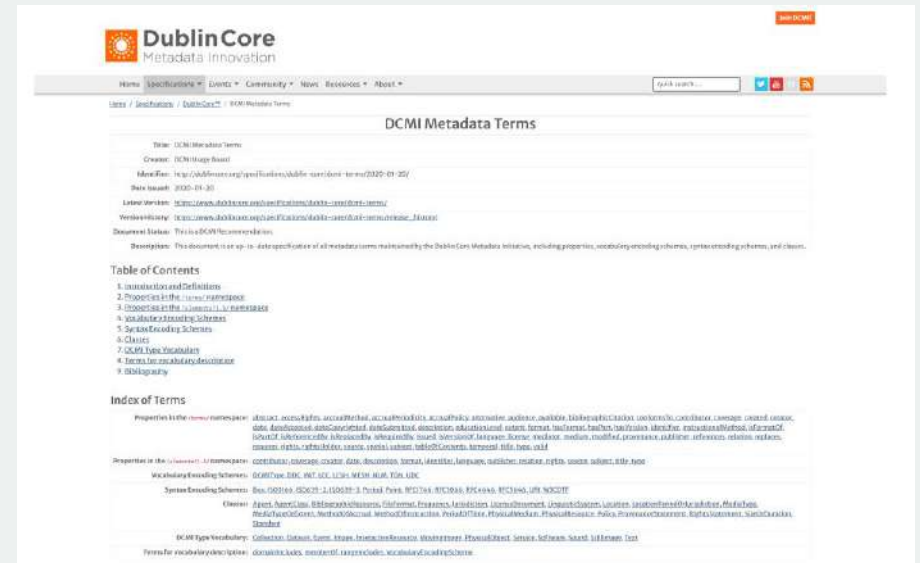
Dublin Core

Core set of element for describing *resources*

- digital resources (video, images, web pages, etc.)
- physical (*books, etc.*)

Examples

- Contributor – "An entity responsible for making contributions to the resource".
- Creator – "An entity primarily responsible for making the resource".
- Date – "A point or period of time associated with an event in the lifecycle of the resource".
- Description – "An account of the resource".
- Format – "The file format, physical medium, or dimensions of the resource".
- Identifier – "An unambiguous reference to the resource within a given context".
- Language – "A language of the resource".
- Type – "The nature or genre of the resource".



Term Name: title		More details
URI	http://purl.org/dc/terms/title	
Label	Title	
Definition	A name given to the resource.	
Type of Term	Property	
Has Range	http://www.w3.org/2000/01/rdf-schema#Literal ← Any string	
Subproperty of	<ul style="list-style-type: none"> • Title (http://purl.org/dc/elements/1.1/title) 	

Term Name: type		More details
URI	http://purl.org/dc/terms/type	
Label	Type	
Definition	The nature or genre of the resource.	
Comment	Recommended practice is to use a controlled vocabulary such as the DCMI Type Vocabulary [DCMI-TYPE]. To describe the file format, physical medium, or dimensions of the resource, use the property Format.	
Type of Term	Property	
Subproperty of	<ul style="list-style-type: none"> • Type (http://purl.org/dc/elements/1.1/type) 	

Dublin Core

Index of Terms

Properties in the /terms/ namespace:	abstract , accessRights , accrualMethod , accrualPeriodicity , accrualPolicy , alternative , audience , available , bibliographicCitation , conformsTo , contributor , coverage , created , creator , date , dateAccepted , dateCopyrighted , dateSubmitted , description , educationLevel , extent , format , hasFormat , hasPart , hasVersion , identifier , instructionalMethod , isFormatOf , isPartOf , isReferencedBy , isReplacedBy , isRequiredBy , issued , isVersionOf , language , license , mediator , medium , modified , provenance , publisher , references , relation , replaces , requires , rights , rightsHolder , source , spatial , subject , tableOfContents , temporal , title , type , valid
Properties in the /elements/1.1/ namespace:	contributor , coverage , creator , date , description , format , identifier , language , publisher , relation , rights , source , subject , title , type
Vocabulary Encoding Schemes:	DCMIType , DDC , IMT , LCC , LCSH , MESH , NLM , TGN , UDC
Syntax Encoding Schemes:	Box , ISO3166 , ISO639-2 , ISO639-3 , Period , Point , RFC1766 , RFC3066 , RFC4646 , RFC5646 , URI , W3CDTF
Classes:	Agent , AgentClass , BibliographicResource , FileFormat , Frequency , Jurisdiction , LicenseDocument , LinguisticSystem , Location , LocationPeriodOrJurisdiction , MediaType , MediaTypeOrExtent , MethodOfAccrual , MethodOfInstruction , PeriodOfTime , PhysicalMedium , PhysicalResource , Policy , ProvenanceStatement , RightsStatement , SizeOrDuration , Standard
DCMI Type Vocabulary:	Collection , Dataset , Event , Image , InteractiveResource , MovingImage , PhysicalObject , Service , Software , Sound , StillImage , Text
Terms for vocabulary description:	domainIncludes , memberOf , rangeIncludes , VocabularyEncodingScheme

Term Name: Dataset

URI <http://purl.org/dc/dcmitype/Dataset>

Label Dataset

Definition Data encoded in a defined structure.

Comment Examples include lists, tables, and databases. A dataset may be useful for direct machine processing.

Type of Term Class

Member Of: <http://purl.org/dc/terms/DCMIType>

DCAT

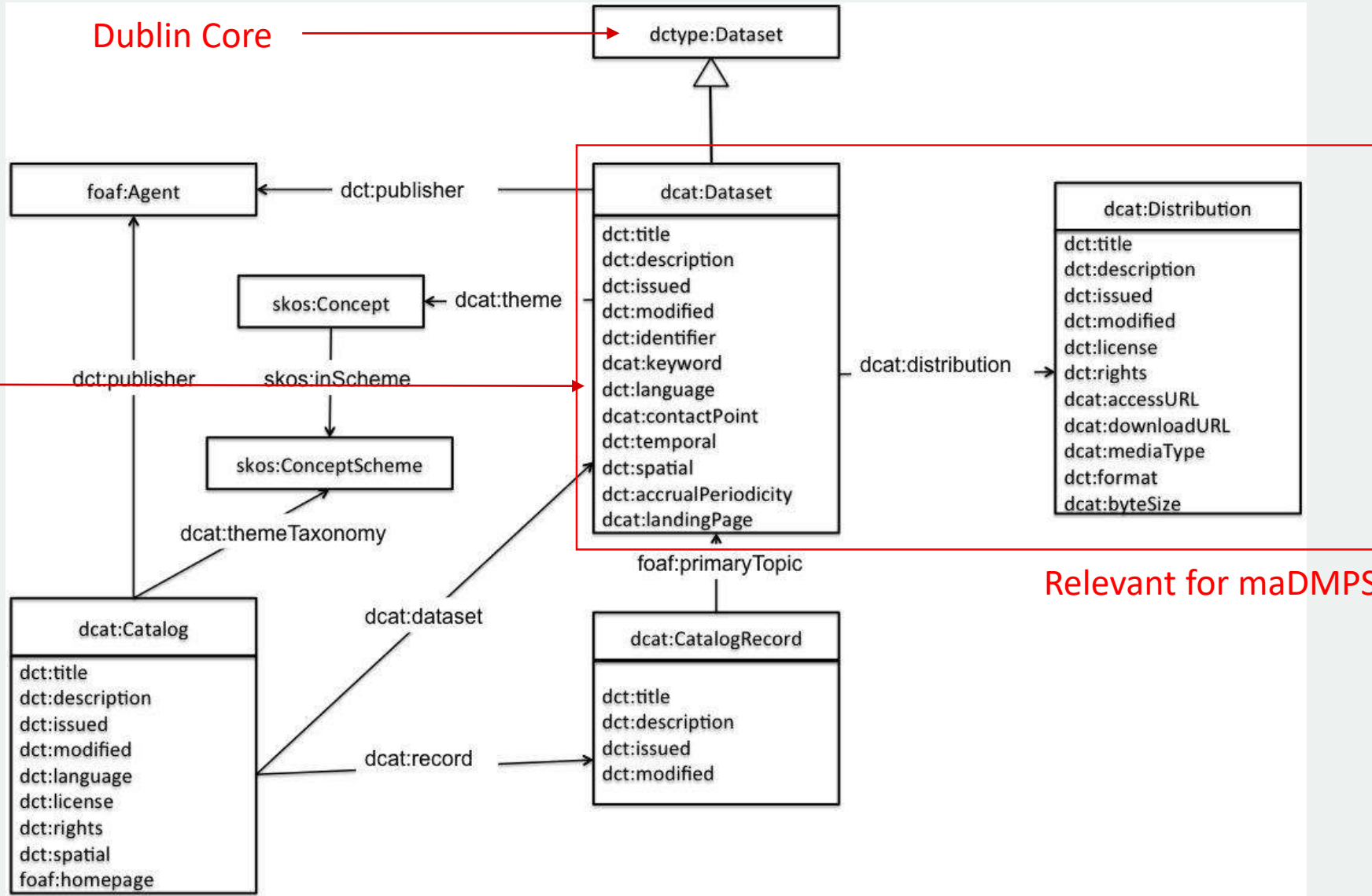
- DCAT is an RDF vocabulary designed to facilitate interoperability between data catalogs published on the Web.
 - Open data portals in Europe did not have the same schema originally
 - <https://data.gov.uk/>
 - Potential use case: support federated search
- No serialization enforced
- Reuses other vocabularies as well, e.g. Dublin Core



A screenshot of the 'Open Data Österreich' website. The page has a blue header with the title 'Open Data Österreich' and a search bar. Below the header, there are three large numbers: 38.345 Datensätze, 688 Anwendungen, and 1.446 Organisationen. The main content area shows a featured dataset 'Bauperiode (Linz)' with a date of 15.02.2022 and a brief description. Below this, there is a section for 'Zufällige Anwendungen basierend auf offenen Daten' which includes a map and a link to a 'Shock' application.

DCAT

Dublin Core



Note the difference:
dcat and *dct*

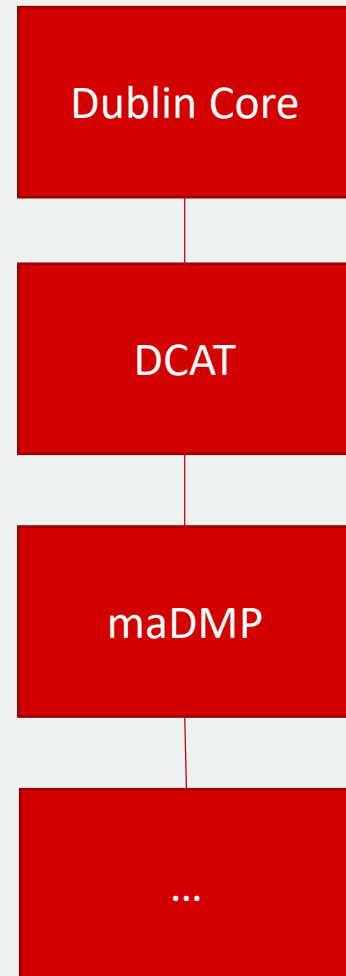
Relevant for maDMPS

<https://www.w3.org/TR/vocab-dcat-1/>

(newer version exists, adds more concepts, does NOT make this one obsolete)

Application Profile

- Adds additional constraints
- Does NOT break compliance
- For example
 - Cardinality constraints (MAY -> MUST)
 - Sub-classes
 - Use further vocabularies
- For this reason DCAT has most of the fields optional
- DCAT itself uses Dublin Core in this way
- maDMP recommendation does the same
 - Reuses classes
 - Sets constraints
 - Adds new terms

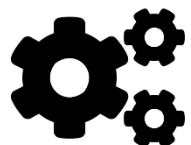


10 PRINCIPLES FOR MADMPS

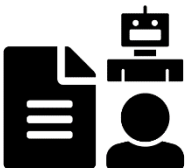
10 principles for maDMPs



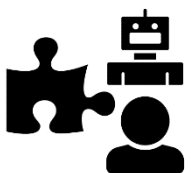
1 Integrate DMPs with the workflows of all stakeholders in the research data ecosystem



2 Allow automated systems to act on behalf of stakeholders



3 Make policies (also) for machines, not just for people



4 Describe—for both machines and humans—the components of the data management ecosystem

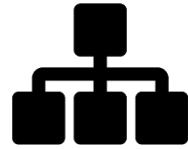


5 Use PIDs and controlled vocabularies

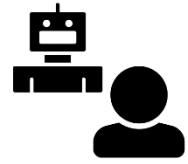
Miksa T, Simms S, Mietchen D, Jones S (2019) Ten principles for machine-actionable data management plans. PLoS Comput Biol 15(3): e1006750.

<https://doi.org/10.1371/journal.pcbi.1006750>

10 principles for maDMPs



6 Follow a common data model for maDMPs



7 Make DMPs available for human and machine consumption



8 Support data management evaluation and monitoring



9 Make DMPs updatable, living, versioned documents



10 Make DMPs publicly available



Miksa T, Simms S, Mitchen D, Jones S (2019) Ten principles for machine-actionable data management plans. PLoS Comput Biol 15(3): e1006750.

<https://doi.org/10.1371/journal.pcbi.1006750>



TECHNISCHE
UNIVERSITÄT
WIEN

FAIR Data Austria

RDM INFRASTRUCTURE

RDM Infrastructure @ TU Wien

DMPs are not for funders only

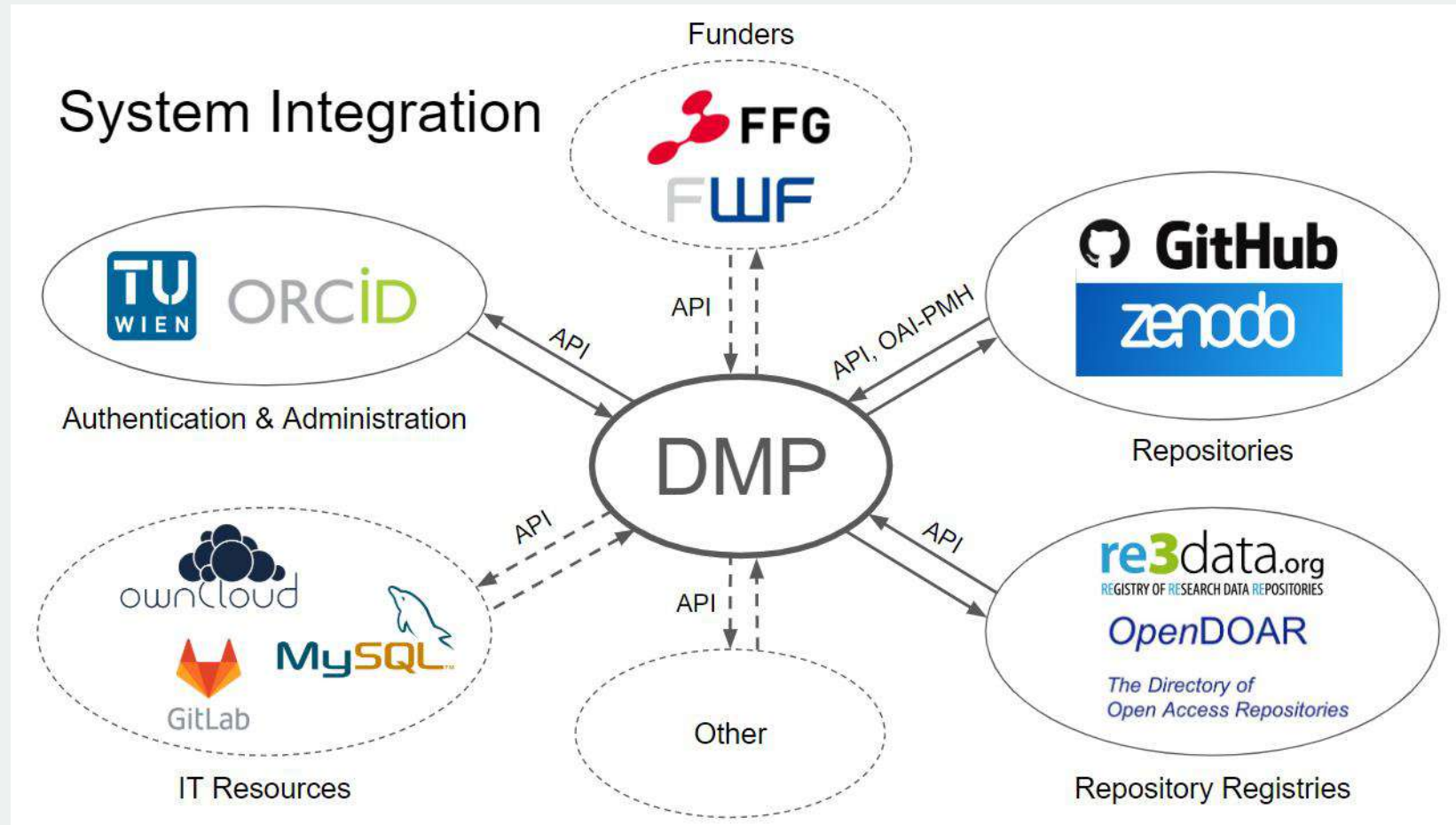
DMPs must also create benefits for researchers

- Less work
- Automation of tasks
- Reuse of information

DMPs are the 'glue' between different systems


- Automate
 - getting data in
 - getting data out

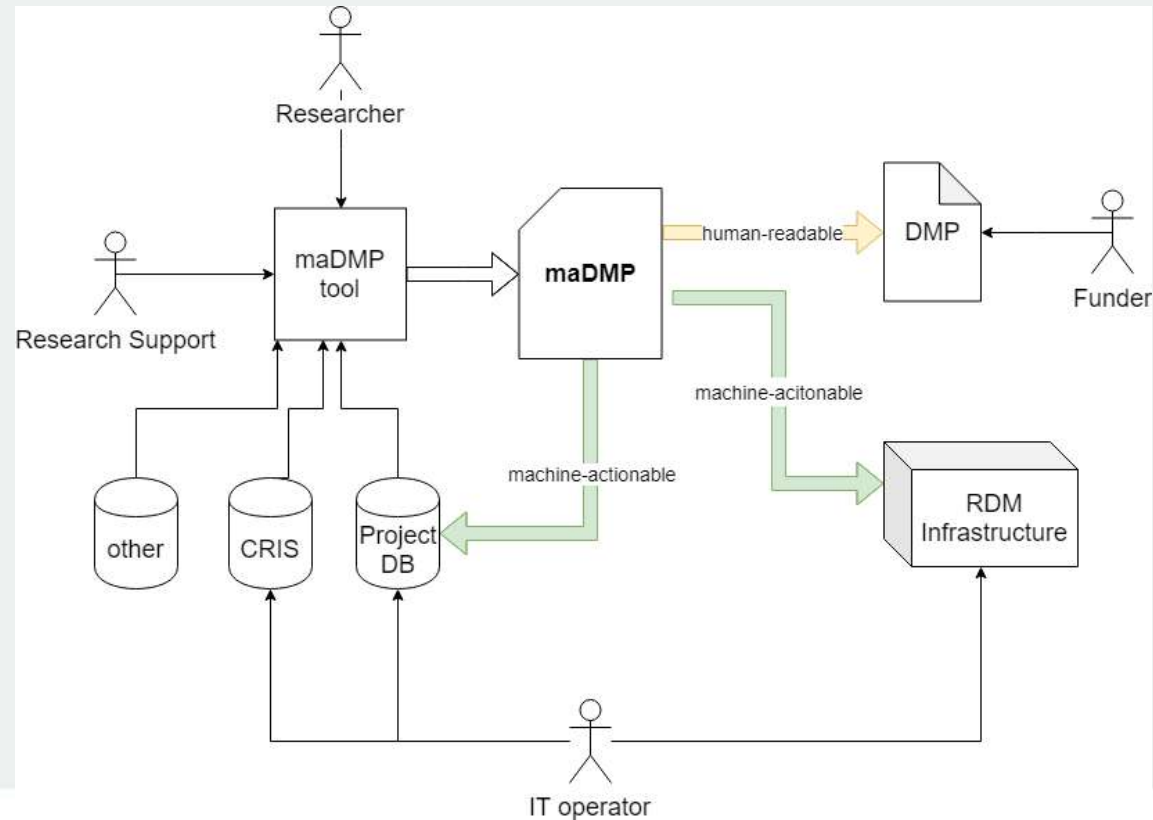
RDM Infrastructure @ TU Wien



maDMPs and RDM infrastructure

once-only principle

- do not ask researchers same questions in different places 



Mock-ups

Requirements collection

- Community
- TU Wien interviews

To be used for MVP

Implemented in DMap

The screenshot displays a web browser window titled "DMP Funder View" with the URL <https://dmpviewer.tuwien.ac.at/dmps/54365437012341>. The page content includes:

- DMP Funder View**: A breadcrumb trail "Home > DMPs > DMP#54365437012341".
- Reuse of pre-existing data**: A table with columns "Dataset title", "Origin", and "License".

Dataset title	Origin	License
Calculating Thermal Bremsstrahlung Emission from Stellar Winds	doi:10.5281/zenodo.1476587	MIT
Occurrence records download on 2018-11-05	doi:10.26197/5be00801ee357	CC-BY
- FAIR Data**:
 - Metadata standards**:
 - [Dublin Core](#)
 - [DataCite Metadata Schema](#)
 - [DDI - Data Documentation Initiative](#)
 - [CIF \(Crystallographic Information Framework\)](#)
 - [CSMD \(Core Scientific Metadata Model\)](#)
 - Metadata**: Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.
- Inferred FAIRness by repository selection**:

Selected repository	Dataset	Data access	PID system	AID system	Certificate	Quality Mgmt.	Versioning	Location	API
GitHub	Source code for client application	open	none	none	none	no	yes	US	other
Zenodo	Supplementary material	open	DOI	ORCID	none	yes	yes	EU	REST OAI-PMH
GESIS Data Archive	Raw data Processed data	open	DOI	none	CoreTrustSeal	no	-	Germany	OAI-PMH
- Licensing**:

Dataset	Sharing strategy	Selected license	License planned to be active from
Supplementary material	keep closed	-	-
Raw data	keep closed	-	-
Source code for client application	publish	Apache License 2	2020-01-01
Processed data	publish	Creative Commons Attribution (CC-BY)	2021-03-01

A yellow callout box in the bottom right corner says "Please click on the scrollbar to see more." with an arrow pointing to the scrollbar.

<https://oblassers.github.io/dmap-mockups/>

DMap

The screenshot displays the DMap web interface. At the top, there is a blue navigation bar with the text "DMap", "My DMP's", and "Create DMP". On the right side of the bar, the user's name "Andreas Rauber" is visible with a dropdown arrow. Below the navigation bar, the main content area shows "You specified 2 dataset(s)".

Two dataset cards are displayed:

- Sumatra Image Collection**: Size range "20GB - 50GB". Description: "Images (JPEG, JPEG2000, GIF, TIF, PNG, SVG, etc.)". Sub-description: "Satellite images of northwestern Sumatra". Size range: "20 - 50 GB".
- Processing App**: Size range "100MB - 1GB". Description: "Source code (scripting, Java, C, C++, Fortran, etc.)". Sub-description: "Python scripts". Size range: "100 - 1000 MB". Description: "Configuration data (parameter settings, logs, library files)". Sub-description: "Configuration files to setup the runtime environment". Size range: "< 100 MB".

Below the dataset cards, there is a section titled "Create names for your datasets" with two input fields: "Sumatra Image Collection" and "Processing App". To the right of these fields is a close button "X".

Below this section, there is a text prompt: "Estimate the type, format and volume of your research data manually and / or by uploading sample data and group them into datasets." Below this prompt are two tabs: "MANUALLY" and "UPLOADING SAMPLE DATA". The "UPLOADING SAMPLE DATA" tab is active, showing a large red dashed border area with a cloud upload icon and the text "Drop files here to analyze".

At the bottom of the interface, there are two buttons: "PREVIOUS STEP" and "NEXT STEP".

At the very bottom, there is a small icon and the text "Documentation and data quality".

<https://www.rd-alliance.org/system/files/documents/2019-RDA-DMAP-Oblasser.pdf>

Test environment

This application instance is for development and testing purposes only. Content may be deleted at any point in time without prior notification.

Welcome to DAMAP, a service that helps you to create and update the Data Management Plan (DMP) for your project.

My DMPs

+ Create new DMP

What is a DMP?

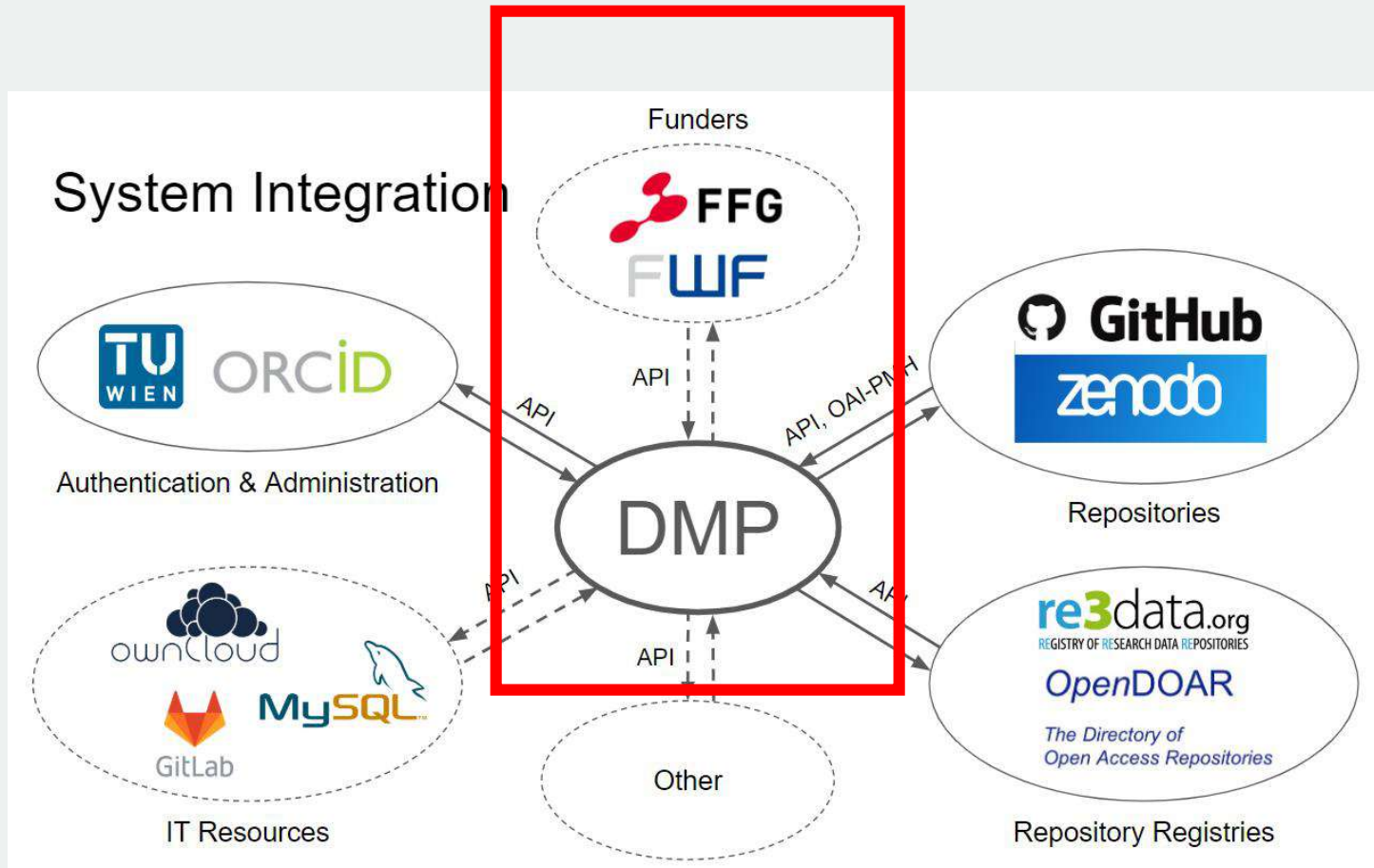
A [DMP](#) is a structured document that keeps record of what research data is created and what happens to that data during and after a project. It helps with planning the research process, managing your data in accordance with the [FAIR Principles](#), and defining rights and responsibilities in a research project involving several researchers or institutions.

DAMAP

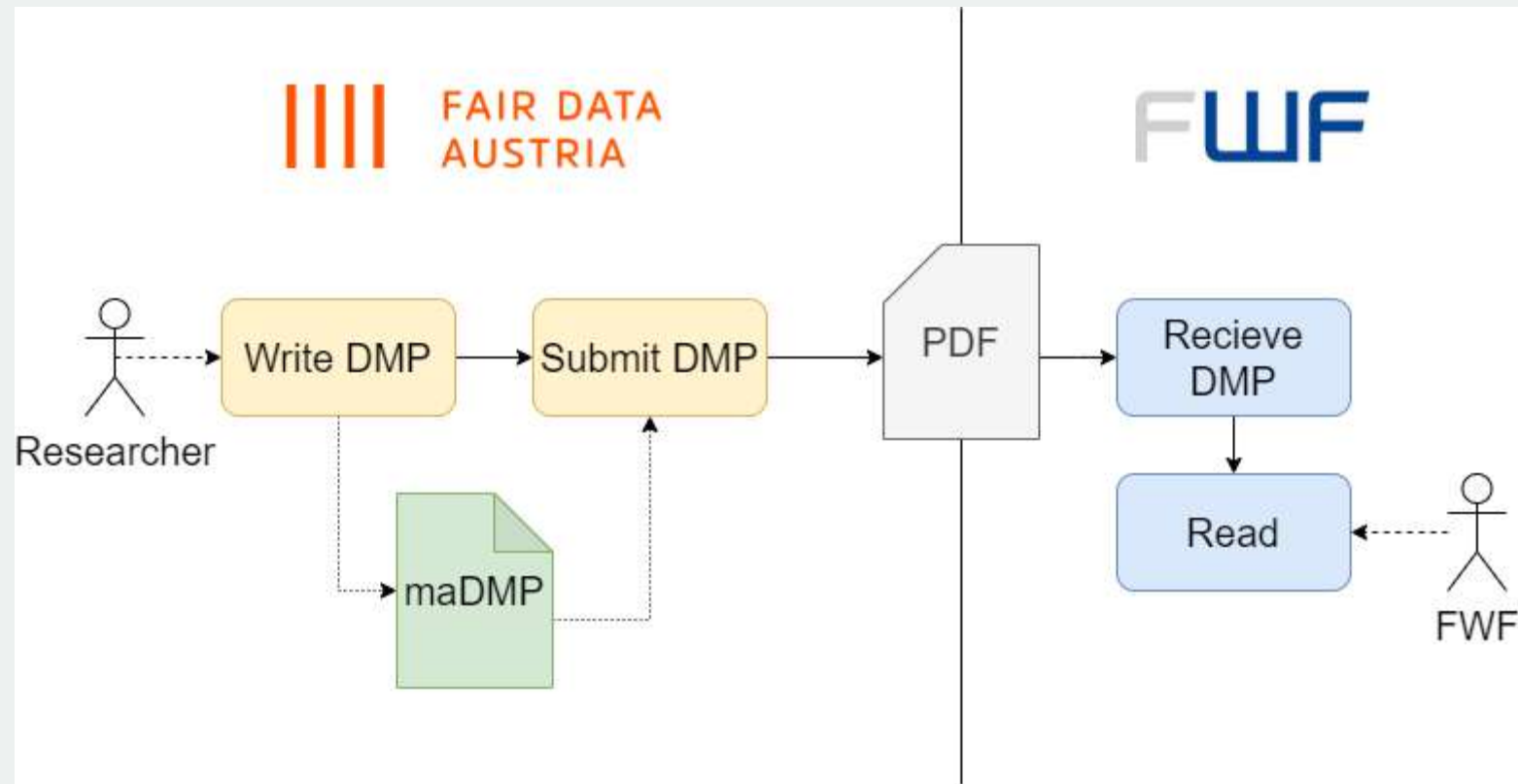
- guides you step by step through the different sections of a DMP following the [Science Europe Practical Guide](#)
- exports a pre-filled DMP as a Word document that you can customize and use for submission to European and national funders, for example FWF and FFG
- saves you work by
 - pre-filling content with detailed information from your CRIS application and other systems
 - providing wizards, guidance, and item lists to choose from
 - suggesting answers that you can either comply with or adjust to your needs
- is compatible with the [RDA recommendation on machine actionable DMPs](#).

<https://damap-qa.apps.dev.csd.tuwien.ac.at>

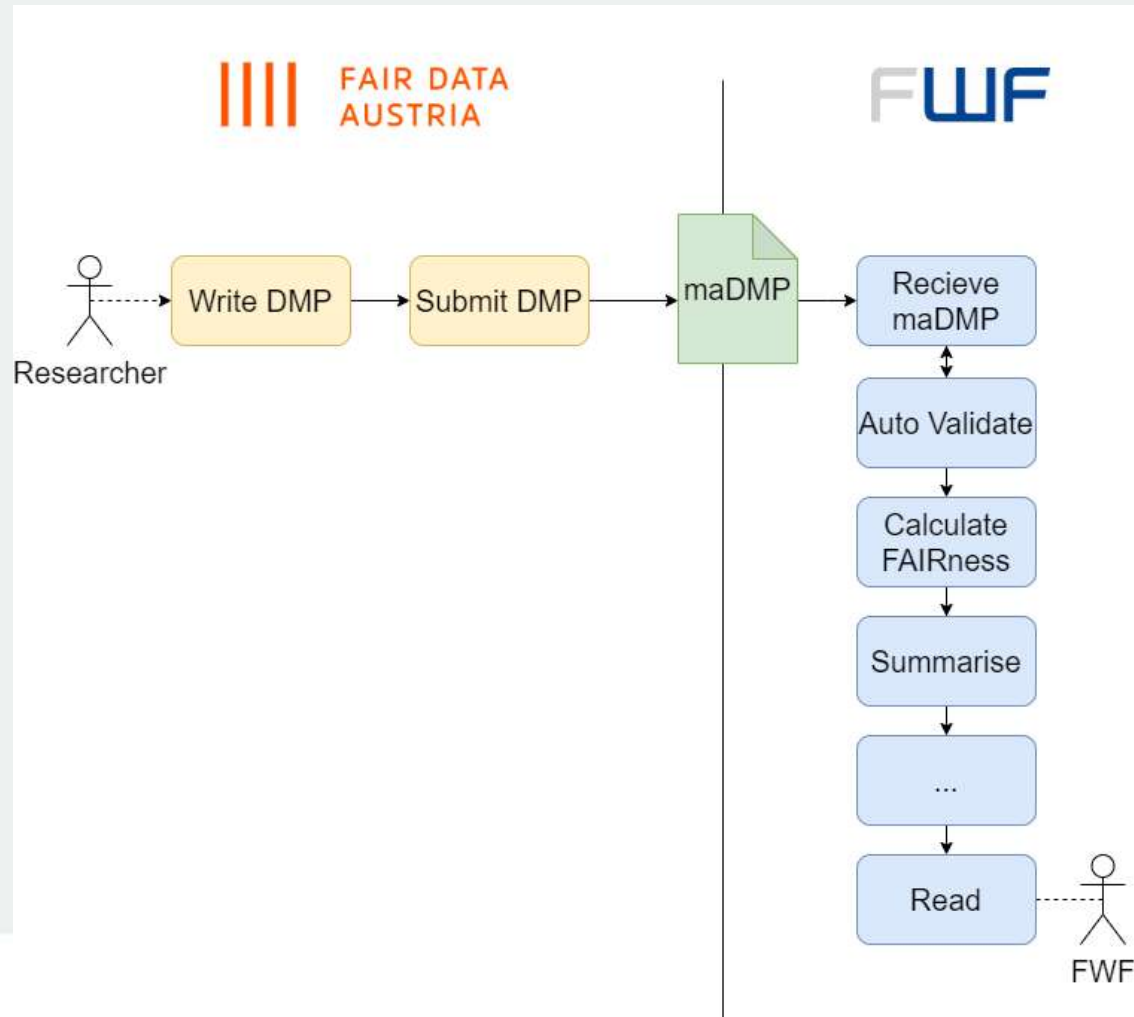
Next steps: Funder integration



Potential scenarios - mixed



Potential scenarios - full



SUMMARY

You should know

Why traditional DMPs are not enough

What is required to make DMPs machine-actionable

What are the related standards and how application profiles work

How common standard for maDMPs works

10 principles for maDMPs

Publications

- [Tomasz Miksa, Simon Oblasser, and Andreas Rauber. **Automating research data management using machine-actionable data management plans**. ACM Transactions on Management Information Systems, 13\(2\), dec 2021.](#)
- [Tomasz Miksa, Paul Walk, Peter Neish, Simon Oblasser, Hollydawn Murray, Tom Renner, Marie-Christine Jacquemot-Perbal, João Cardoso, Trond Kvamme, Maria Praetzellis, Marek Suchánek, Rob Hooft, Benjamin Faure, Hanne Moa, Adil Hasan, and Sarah Jones. **Application profile for machine-actionable data management plans**. CODATA Data Science Journal, 20\(1\):32, October 2021](#)
- [João Cardoso, Leyla Jael Castro, and Tomasz Miksa. **Interconnecting Systems Using Machine-Actionable Data Management Plans Hackathon Report**. Data Science Journal, 20, 2021.](#)
- [Tomasz Miksa, Maroua Jaoua, and Ghaith Arfaoui. **Research Object Crates and Machine-actionable Data Management Plans**. In DaMaLOS - First Workshop on Data and Research Objects Management for Linked Open Science : Co-located at the International Semantic Web Conference ISWC 2020. PUBLISSO, November 2020.](#)
- [Simon Oblasser, Tomasz Miksa, Asanobu Kitamoto: **Finding a Repository with the Help of Machine-Actionable DMPs: Opportunities and Challenges**. IDCC 2020](#)
- [Tomasz Miksa, Stephanie Simms, Daniel Mietchen, Sarah Jones \(2019\) **Ten principles for machine-actionable data management plans**. PLOS Computational Biology 15\(3\): e1006750.](#)
- [Tomasz Miksa, Peter Neish, Paul Walk, Andreas Rauber: **Defining requirements for machine-actionable Data Management Plans**. iPres 2018](#)
- [Tomasz Miksa, João Cardoso, José Luis Borbinha: **Framing the scope of the common data model for machine-actionable Data Management Plans**. BigData 2018: 2733-2742](#)

A Database Repository for Sensitive Data

Martin Weise

Information and Software Engineering
Technische Universität Wien

May 9th, 2022

Introduction

About me

This is my first lecture!

- ▶ PreDoc Researcher @IFS
- ▶ MSc in *Software Engineering & Internet Computing* 2022 (TU Wien)
- ▶ BSc in *Software & Information Engineering* 2019 (TU Wien)



Digital repositories frequently struggle to make **databases** available in their collection:

1. **Separation of concerns:** data stewards receive a database without semantic understanding of the data, curation activities start already poorly
2. **Project phases:** deployed at research unit level, maintained by researchers and handover to IT staff once finished
3. **FAIR queries:** preserving databases, FAIR principles become a challenge
4. **Reproducibility:** re-execution of a (persistently identified) query to reproduce research output is often times neglected, preservation is used in a archiving context
5. **Versioning:** monolithic records, large administrative overhead through lack of machine-readable interfaces

Digital repositories also struggle to make **sensitive data** available in their collection:

1. **Control:** maintaining protection over sensitive data
2. **Open science:** requires data sharing, conflicts with allowing access to third parties
3. **Data sharing:** give up control of the data, typically once this happens the data is gone

Definition: Sensitive Data [1]

Any ideas?

Digital repositories also struggle to make **sensitive data** available in their collection:

1. **Control:** maintaining protection over sensitive data
2. **Open science:** requires data sharing, conflicts with allowing access to third parties
3. **Data sharing:** give up control of the data, typically once this happens the data is gone

Definition: Sensitive Data [1]

“[...] any information that is protected against unwarranted disclosure. Protection of data may be required for legal or ethical reasons, for issues pertaining to personal privacy, or for proprietary considerations.”

Introduction

What is Data Visiting?

- ▶ Data stays under the control of the owner
 - Who can access
 - Over which period
 - Which subset of data
 - Answer which research question(s) and activities
- ▶ Allow consumers (i.e. analysts, machine learning algorithms) to come to the data
 - Ensure proper usage of data (=minimize risk)
 - Allow **maximum usage of data**
- ▶ Closely monitoring processes and interaction with the data during these visits
- ▶ Safe-guards to prevent accidental leakage and **intentional data breaches**

Introduction

Example: Intentional Data Breaches

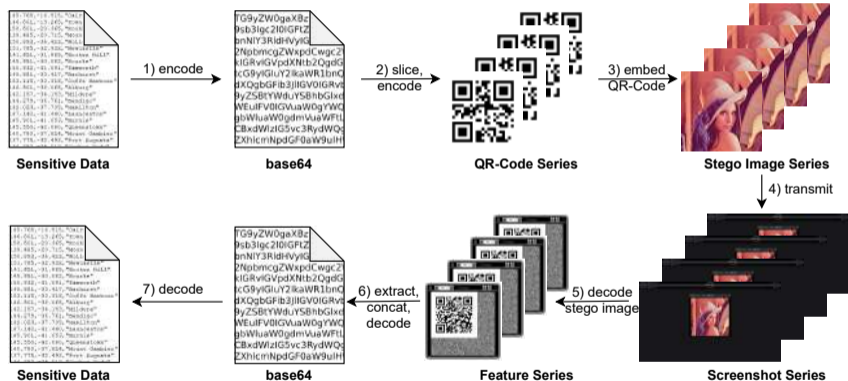


Figure: QR-Code Optical Covert Channel [2]

(Overview)

- ▶ *Infrastructure*: Openstack, Docker
- ▶ *Databases*: MariaDB, Postgres
- ▶ *Automatization*: Ansible, Python, Bash
- ▶ *Identity management*: FreeIPA
- ▶ *Retrieval*: ElasticSearch
- ▶ *Data visiting*: OpenVPN, TigerVNC

What makes this a repository?

We use a similar stack as most web shops nowadays, how does this differ?

Introduction

What makes this a repository

Technical Measures

Safeguarding data against technical threats and design flaws, provide a highly-controlled virtual research environment

Introduction

What makes this a repository

Technical Measures

Safeguarding data against technical threats and design flaws, provide a highly-controlled virtual research environment

Organizational Measures

Identification of researchers, research questions and guide all involved people through well-defined processes on need-to-know basis

Introduction

What makes this a repository

Technical Measures

Safeguarding data against technical threats and design flaws, provide a highly-controlled virtual research environment

Organizational Measures

Identification of researchers, research questions and guide all involved people through well-defined processes on need-to-know basis

Legal Measures

Build trust with data providers, provide legal framework to discourage information sharing, provide foundation for legal steps after (accidental) leaks

Technical Measures

Example: Homomorphic Encryption

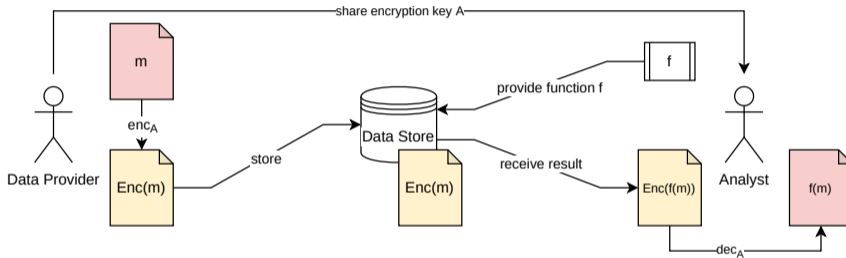


Figure: Simplified homomorphic encryption scheme with a *untrusted* Data Store

The Data Store is now hidden from the Analyst

Is this sufficient for most researchers?

Technical Measures

Example: Homomorphic Encryption

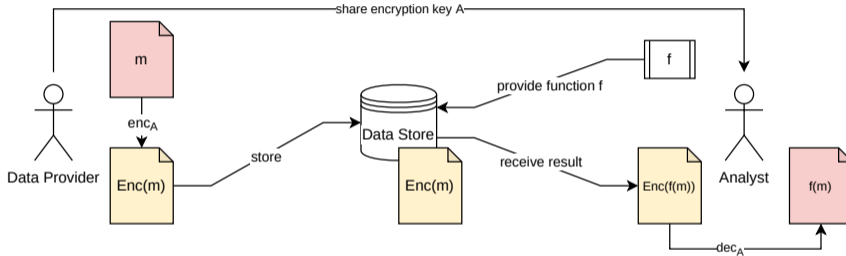


Figure: Simplified homomorphic encryption scheme with a *untrusted* Data Store

The Data Store is now hidden from the Analyst

Still not sufficient for exploratory/interactive analyses types required in many settings!

Technical Measures

Example: Secure Enclaves

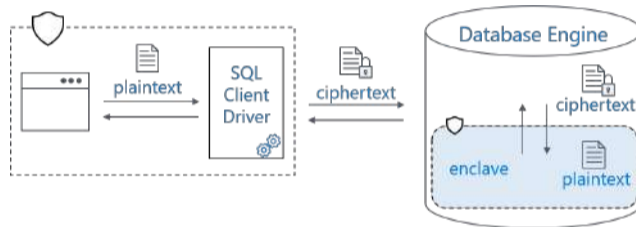


Figure: Permanent encryption of sensitive data in Azure SQL Database [3]

Secure enclave computation on untrusted hardware

Provide hardware isolation and memory encryption, isolating application code and data from untrusted hardware. Many manufacturers of CPUs offer secure enclave support for consumer electronics.

“Five safe” dimensions [4] to maximize public value while also protecting personal rights of individuals and allowing flexible solutions:

- ▶ *Safe projects*: management decisions regarding appropriateness of the usage of the data through auditability and review processes.
- ▶ *Safe people*: identify individuals that access the sensitive data and require them to sign legally binding terms of use.
- ▶ *Safe data*: ensure appropriate data de-identification and access capabilities with respect to the research questions formulated.
- ▶ *Safe settings*: address the necessity of security and transparency to achieve trust with the public and data owners.
- ▶ *Safe outputs*: ensure only approved, aggregated research results can be exported.

Technical Measures

Data Visiting: Air-Gap Isolation



Figure: Schematic overview example of a TRE environment [5]

Technical Measures

Data Visiting: Air-Gap Isolation

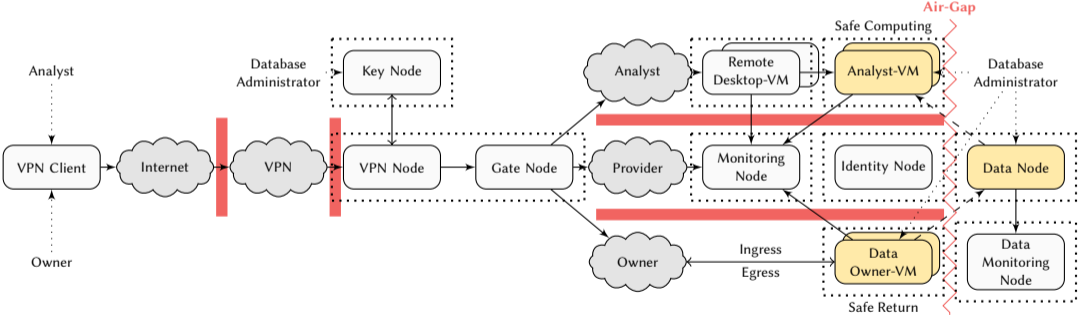


Figure: Our open-source technical architecture specification of OSSDIP [6]

Technical Measures

Data Versioning

ID	Sensor	Temp
1	A	23.1
2	B	25.8

(a) Original Table

ID	Sensor	Temp
1	A	22.1
2	B	25.8

(b) Corrected Table

Figure: Update operation on a table, creating a new version at a different timestamp

Why do we need data versioning? Why not use snapshots to cite them?

Any ideas?

ID	Sensor	Temp
1	A	23.1
2	B	25.8

(a) Original Table

ID	Sensor	Temp
1	A	22.1
2	B	25.8

(b) Corrected Table

Figure: Update operation on a table, creating a new version at a different timestamp

Why do we need data versioning? Why not use snapshots to cite them?

Full copies pollute the database, selection of the “right” snapshot difficult knowing only the query, challenges to archive every single snapshot, ...

Technical Measures

Data Versioning (cont.)

ID	Sensor	Temp	Valid From	Valid To
1	A	23.1	t1	
2	B	25.8	t2	

(a) Original Table

ID	Sensor	Temp	Valid From	Valid To
1	A	23.1	t1	t3
2	B	25.8	t2	
1	A	22.1	t3	

(b) Corrected Table

Figure: Implementation of data versioning in MariaDB (v10.5)

Machine-readability and -actionability through APIs:

- ▶ *Hypertext Transfer Protocol (HTTP) API*
 - REST constraints: client-server architecture, stateless, uniform resource identification
 - Service discovery: 8 services and 14 endpoints visible in the gateway
 - Allow any researcher to (re-)execute view-only queries [7]
- ▶ *Advanced Message Queue Protocol (AMQP) API*
 - Time-series tuple insert (i.e. from sensors)
 - One exchange per database, one queue per table, handover to HTTP API
- ▶ *Java Database Connectivity (JDBC) API*
 - Direct access for database experts to maintain the database

Interaction with external systems:

- ▶ DOI: DataCite schema to persistently identify a query and *mint a DOI* (soon)

Organizational Measures

Ingress Sensitive Data into the System

Data Owner wants to import a dataset into the infrastructure:

- ▶ Must sign a *data processing agreement*, receive training material, etc.
- ▶ Provide dataset (and metadata) via the **GUI**, receive an account
- ▶ Deposit into Data Owner-VM
- ▶ Database Administrator briefly connects Data Node and copies dataset, restores air-gap

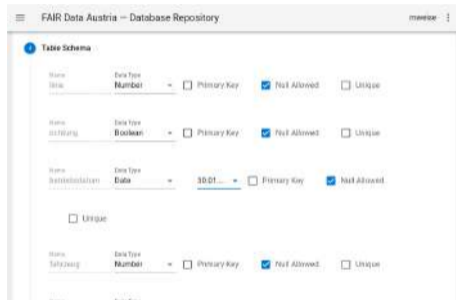


Figure: Provide dataset metadata

Organizational Measures

Ingress Sensitive Data into the System

Add metadata to the table schema:

- ▶ For now, only units of measurement [8] (centimeter, kilogramm, etc.)
- ▶ Assign a **unit to each column** (where applicable)

Why is this important?

Any ideas?

Column Unit for "kurs"

degree

Name
degree

Symbol
-

Comment
The degree is a unit of angle defined as 1.745329e-2 radian.

URI
<http://www.ontology-of-units-of-measure.org/resource/om-2/degree>

CLOSE SAVE

Figure: Assign a unit of measurement

Organizational Measures

Ingress Sensitive Data into the System

Add metadata to the table schema:

- ▶ For now, only units of measurement [8] (centimeter, kilogramm, etc.)
- ▶ Assign a **unit to each column** (where applicable)

Why is this important?

We assign specific semantic, machine-readable concept to a column that allows conversion to other concepts, e.g.
 $100 [^{\circ}\text{C}] = 212 [^{\circ}\text{F}]$



Figure: Assign a unit of measurement

Organizational Measures

Ingress Sensitive Data into the System

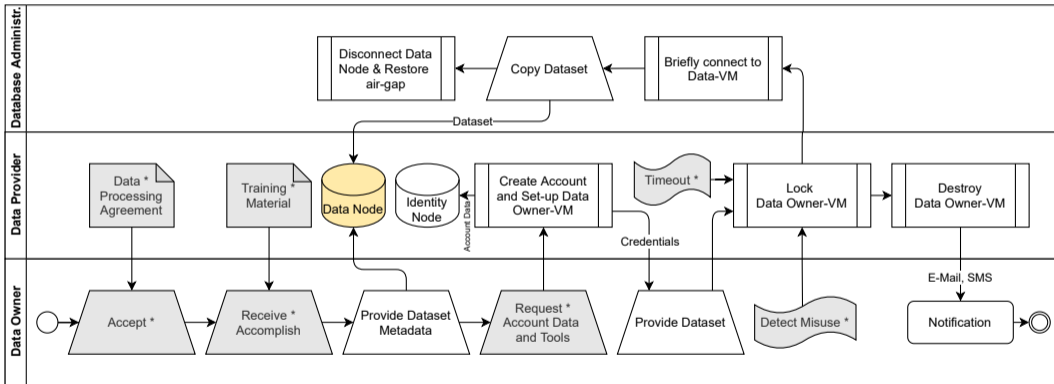


Figure: Complete data ingress process [6]

Organizational Measures

Requesting Access to Sensitive Data

1. Findability

- F1: Assign a PID to the query that creates a subset (for the Analyst)
- F2: Add unit of measurement, require metadata on data upfront
- F3: Link PID of data along with the metadata
- F4: Metadata is searchable by ElasticSearch

2. Accessibility

- A1: open and authentication via HTTPS, AMQP, JDBC
- A2: metadata always available in metadata database

3. Interoperability

- I1: OWL/RDF concepts for units of measurement
- I2: no metadata on the services yet, in development [9]
- I3: metadata is interlinked in the metadata database

4. Reuseability

- R1: open-source license, describe the data (currently only free form)

FAIR Sensitive Data

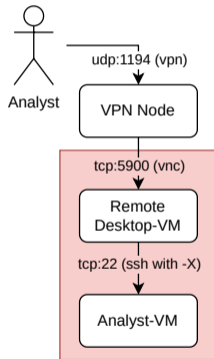
Even though the data is not open, the dataset is still FAIR!

Organizational Measures

Requesting Access to Sensitive Data

Analyst wants to access a sensitive dataset:

- ▶ Send personal identification data along with proposal
 - Required sensitive data (metadata is known)
 - Required tools to analyze data
- ▶ Extract the required subset via query (and store this query in the query store with timestamp, etc.)
- ▶ Potentially apply fingerprinting to the subset [10]
- ▶ Provision of Analyst-VM with the imported subset and Remote Desktop-VM



Organizational Measures

Requesting Access to Sensitive Data

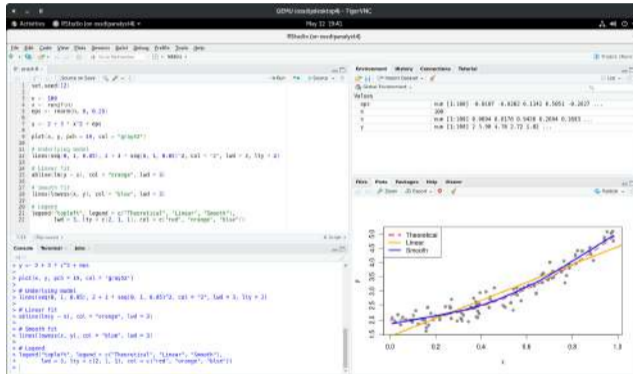


Figure: Work with the data with pre-approved tools via VPN+VNC

Organizational Measures

Requesting Access to Sensitive Data

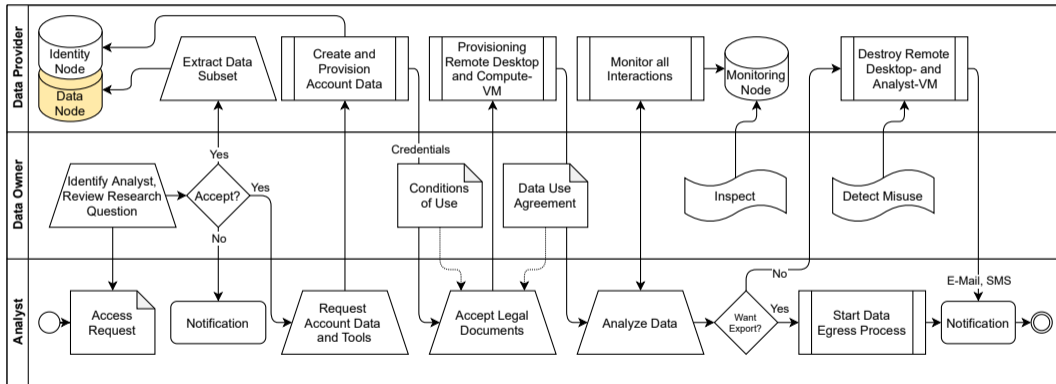


Figure: Complete data access process [6]

Organizational Measures

Subset Generation

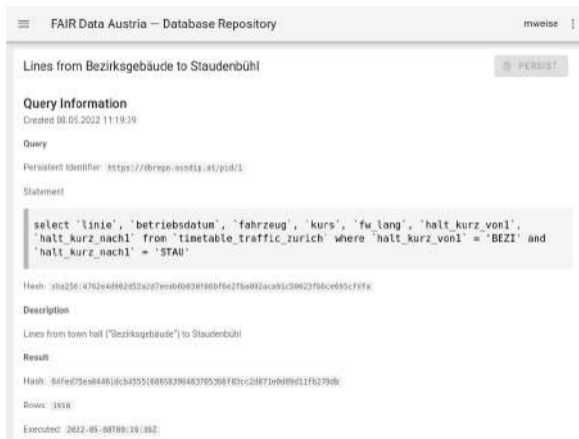
The screenshot displays the 'FAIR Data Austria – Database Repository' interface. At the top, there is a navigation menu and a user profile 'mweide'. The main section is titled 'Create Query' and features a blue 'EXECUTE' button. Below this, a table named 'Timetable Traffic Zürich' is selected, showing columns: 'linie', 'betriebsdatum', 'fahrzeug', 'kurs', 'fw_lang', 'halt_kurz_von1', and 'halt_kurz_nach1'. The query is constructed using a visual builder with two 'where' clauses: 'halt_kurz_von1 = BEZI' and 'halt_kurz_nach1 = STAU', connected by an 'and' operator. The resulting SQL query is shown in a text area below:

```
select
  'linie'
  'betriebsdatum'
  'fahrzeug'
  'kurs'
  'fw_lang'
  'halt_kurz_von1'
  'halt_kurz_nach1'
from
  'timetable_traffic_zurich'
where
  'halt_kurz_von1' = 'BEZI'
and
  'halt_kurz_nach1' = 'STAU'
```

Figure: Lines from town hall (“Bezirksgebäude”) to Staudenbühl

Organizational Measures

Persistent Identification of Subsets



FAIR Data Austria – Database Repository mweise

Lines from Bezirksgebäude to Staudenbühl PERIST

Query Information
Created 08.05.2022 11:19:29

Query
Persistent Identifier: <https://dbrepo.ozdip.at/pid/1>

Statement

```
select `linie`, `betriebsdatum`, `fahrzeug`, `kurs`, `fw_lang`, `halt_kurz_von1`,  
`halt_kurz_nach1` from `timetable_traffic_zurich` where `halt_kurz_von1` = 'BEZI' and  
`halt_kurz_nach1` = 'STAU'
```

Hash: sha256:4762e4d942d52a297e5b6e03e186bf5e2fba892aca81c58923f3bce685cf7fa

Description
Lines from town hall ("Bezirksgebäude") to Staudenbühl

Result
Hash: bf4ee75ea84461dcb4555168655298483765268f83cc2d871e9d89d11f6279db
Rows: 1818
Executed: 2022-05-08T09:39:19Z

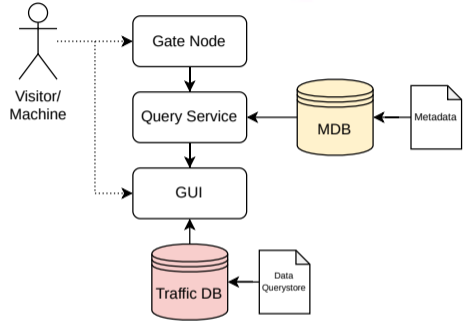
Figure: Query from the query store, persisted in the metadata database

Organizational Measures

Persistent Identification of Subsets

Stored query:

- ▶ Query (“raw” and normalized)
- ▶ Timestamp of creation
- ▶ Timestamp of execution
- ▶ Result hash
- ▶ Creator



How does a stored query differ from a persisted query?

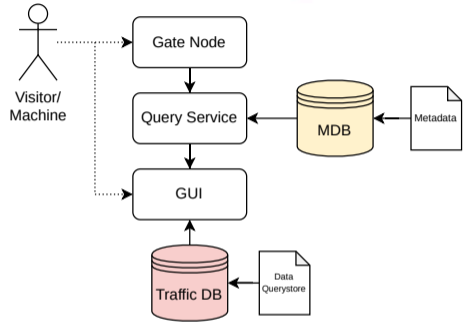
Any thoughts?

Organizational Measures

Persistent Identification of Subsets

Stored query:

- ▶ Query (“raw” and normalized)
- ▶ Timestamp of creation
- ▶ Timestamp of execution
- ▶ Result hash
- ▶ Creator



How does a stored query differ from a persisted query?

Persisted queries store metadata also in the metadata database. When the data is not available anymore, the metadata is still available!

Overall goal: prevent (unauthorized) leaks of sensitive data, what do we need legally?

- ▶ Trusted computing infrastructure for Data Owners
- ▶ Conditions of use for Analysts
 - Prohibition of data download
 - Prohibition of de-anonymization
 - Non-Disclosure Agreement (NDA)
 - Agreement to extensive monitoring
- ▶ Information in case of unauthorized leaks
 - Personal identifiable information
 - Evidence of misuse (compare to agreed research questions)
 - Evidence of relationship between leaked dataset and provided dataset (i.e. through fingerprinting)

	Traditional Repository¹	Database Repository
Representation	File storage	Database engine
Versioning	Snapshots	Temporal tables
Identification	Record	Query/Subset
Storage	Filesystem	Database
Concerns	Researcher	Database expert
Use-case	Deposit	Continuous work

Table: Key differences between a traditional repository a database repositior

¹e.g. Invenio, Dataverse, Figshare, Mendeley Data, Open Science Framework

Conclusion

A Database Repository for Sensitive Data

Important concepts for the exam:

1. Sensitive data and data visiting
2. Three aspects of sensitive data repositories
3. Data versioning
4. Stored and persisted queries
5. FAIR sensitive data

Conclusion

Learn more about these projects:

- ▶ Database Repository
<https://dbrepo-docs.ossdip.at>
- ▶ Secure Data Infrastructure / Data Visiting
<https://ossdip.at>

Talk to me after the lecture or drop me a mail martin.weise@tuwien.ac.at!

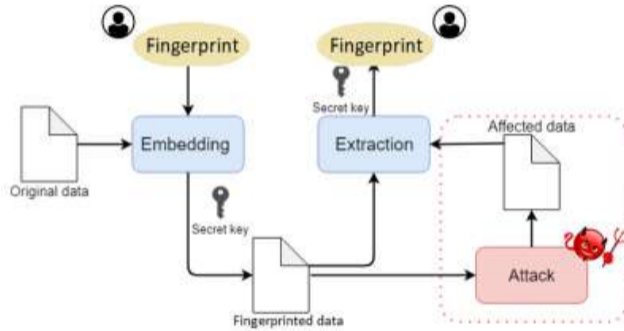


Figure: General fingerprinting scheme [10]

Backup Slides

Marking Ownership (cont.)

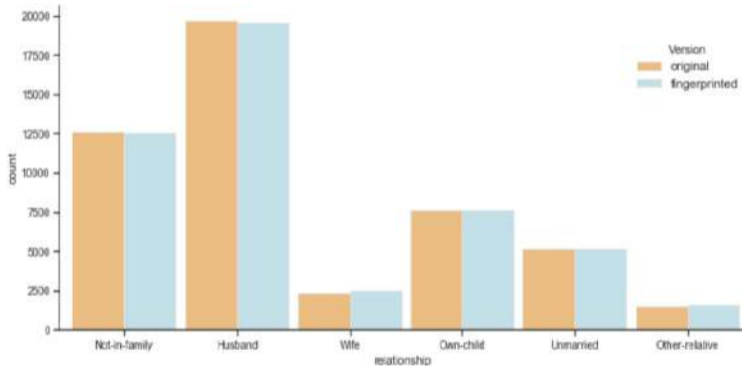










Figure: Distribution of a categorical attribute before and after fingerprinting [10]

-  CSC - IT Center for Science, “Definition of sensitive data,” [Online]. URL: <https://research.csc.fi/definition-of-sensitive-data>, accessed 2022-05-07.
-  M. Weise, “A QR-Code Optical Covert Channel in an Air-Gapped Secure Data Infrastructure,” Master’s thesis, Technische Universität Wien, 2022.
-  J. Szymaszek, “Always Encrypted with Secure Enclaves in Azure SQL Database Preview,” [Online]. URL: <https://techcommunity.microsoft.com/t5/azure-sql-blog/always-encrypted-with-secure-enclaves-in-azure-sql-database/ba-p/2051544>, 2021, accessed 2022-05-07.
-  T. Desai, F. Ritchie, and R. Welpton, “Five Safes: Designing data access for research,” 2016. doi: 10.13140/RG.2.1.3661.1604 Economics Working Paper Series 1601.

-  UK Health Data Research Alliance and NHSX, “Building Trusted Research Environments - Principles and Best Practices; Towards TRE Ecosystems,” 2021. doi: 10.5281/zenodo.5767586
-  M. Weise, F. Kovacevic, N. Popper, and A. Rauber, “Open Source Secure Data Infrastructure and Processes Supporting Data Visiting,” *Data Science Journal*, vol. 21, no. 1, p. 4, 2022. doi: 10.5334/dsj-2022-004
-  A. Trisovic, P. Durbin, T. Schlatter *et al.*, “Advancing Computational Reproducibility in the Dataverse Data Repository Platform,” in *Proceedings of the 3rd International Workshop on Practical Reproducible Evaluation of Computer Systems*, New York, NY, USA, 2020. doi: 10.1145/3391800.3398173 p. 15–20.
-  H. Rijgersberg, “Ontology of Units of Measure,” <https://github.com/HajoRijgersberg/OM>, 2022, accessed 2022-05-08.

Bibliography III

-  L. O. Bonino, K. Burger, and R. Kaliyaperumal, “FAIR Data Point,” [Online]. URL: <https://specs.fairdatapoint.org/#repositorymetadata>, 2021, accessed 2022-05-08, version 1.0.
-  T. Sarcevic and R. Mayer, “Fingerprinting Relational Data Sets,” [Poster]. [Online]. URL: https://www.sba-research.org/wp-content/uploads/2020/09/MLDM_Fingerprinting-Relational-Data-Sets_WEB.pdf, 2020, accessed 2022-05-07.

REPOSITORIES I: TV DATA⁽¹⁾ AND ITS DESIGN

INVENIQ RDM

⁽¹⁾ The name finding process is still ongoing... Let's hope it won't be TURD again

About Me

Maximilian Moser

- From the Tyrolean Alps
- Studied Software Engineering at TU Wien
- Not a lecturer, but a DevOp for TU Data



mmoser@tuwien ~ \$



What is a Research Data Repository?

- Central place for researchers to deposit and share their data
 - No more Google Drives and institute servers...
- For *all kinds* of research data
 - Stores datasets along with *metadata*
- Helps with FAIR-ification of the datasets
 - Provides visibility and access restrictions
- Ensures that data stays available
 - Archival and preservation

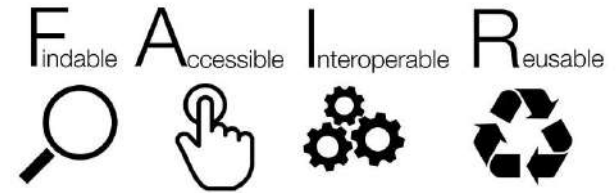
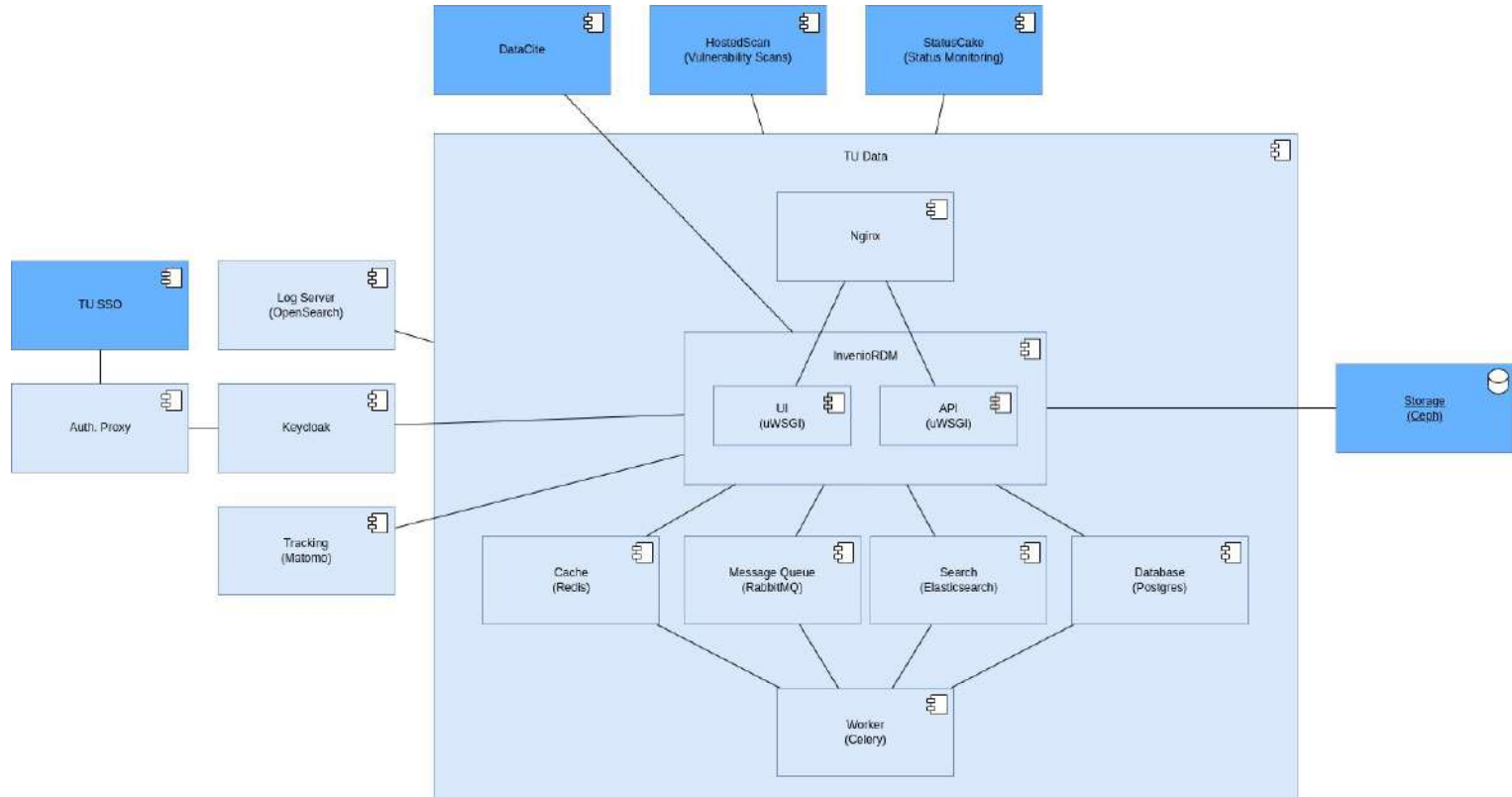


Image by [SangyaPundir](#)

What is InvenioRDM / TU Data?

- File-based research data repository
 - Developed by CERN and partners: <https://inveniosoftware.org/products/rdm/>
- Flask Application
 - Frontend: React, Semantic UI
 - Backend: Python
- Highly customizable and extensible
 - TU Data is a [themed and customized](#) variant of InvenioRDM
- Still under development

Architecture Overview



Frontpage

HOME MY DASHBOARD SETTINGS TU WIEN

Welcome to TU Data Repository

TU Data Repository is an institutional service of TU Wien to enable storing, sharing and publishing of digital objects, in particular research data. It facilitates the broader requirements for open access in research data and the FAIR principles by making research output reusable, accessible, interoperable and so on. This service is developed by the TU Wien Center for Research Data Management and is hosted by TUG.

Please note that this service is still under development and may be subject to change. We will address functionality as development progresses.

Deposit Search FAQ Contact

Recent Uploads

ORCADs
Erich, Wolfgang (Researcher); Janda, Jiri (TU Wien, Austria)
[View] [Download] [Share] [Like]

ORCADs is an accelerated version of ORCADs Dataset (Eckardt et al., 2009) accelerated with user advice using weak supervision. It allows you to test your algorithm on various types of data records. These records are already taken from Berlin's (2002) moviebusan advertisement, (re)generated...

ORCADs Dataset (Eckardt et al., 2009)
uploaded on March 29, 2023

FAIR for Sensitive Data
Meyer, David; Tomczak, Piotr; Chapman, Peter J.; Wambler, Laura; Chalkin, Andrew; Miles, Thomas; Szepietowski, Tomasz
[View] [Download] [Share] [Like]

Meyer's presentation is a set of slides of the earlier website FAIR for Sensitive Data, registered by the FAIR Commons Austria on March 29, 2023. The goal of the website was to inform researchers on technical and legal aspects and best practices when working with sensitive data. The website...

FAIR for Sensitive Data (Eckardt et al., 2009)
uploaded on March 29, 2023

RT-Percept Gas Template
Gomboc, Jan (TU Wien)
[View] [Download] [Share] [Like]

Pre-rendered dataset used in Training and Predicting Visual Error for Real-Time Applications for the Gas Template scene. Generated using the RT-Percept renderer and the RT-Percept system.
uploaded on March 29, 2023

RT-Percept Slovak Cathedral
Gomboc, Jan (TU Wien)
[View] [Download] [Share] [Like]

Pre-rendered dataset used in Training and Predicting Visual Error for Real-Time Applications for the Slovak Cathedral scene. Generated using the RT-Percept renderer and the RT-Percept system.
uploaded on March 29, 2023

RT-Percept Looswerynk-Blaas
Gomboc, Jan (TU Wien)
[View] [Download] [Share] [Like]

Pre-rendered dataset used in Training and Predicting Visual Error for Real-Time Applications for the Looswerynk-Blaas scene. Generated using the RT-Percept renderer and the RT-Percept system.
uploaded on March 29, 2023

INVENTORUM FAIR DATA ALLIANCE POLICIES TERMS OF USE DATA PROTECTION DECLARATION CONTACT

<https://researchdata.tuwien.ac.at>

Search

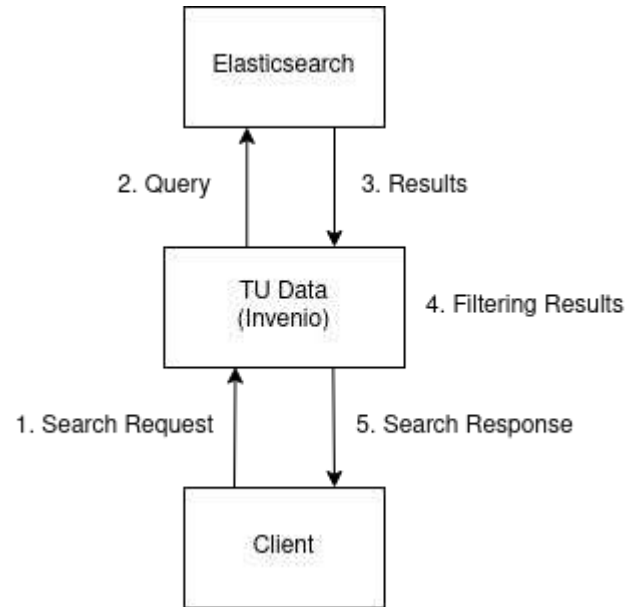


The screenshot displays the Invenio search results page for the query 'CLEF-IP'. The interface includes a search bar at the top with the query 'CLEF-IP' and a search button. Below the search bar, there are filters for 'Versions', 'Access status', 'Resource types', and 'Help'. The search results are listed in a vertical column, each entry featuring a title, author information, a brief description, and a 'Best match' dropdown menu. The results include:

- The CLEF-IP 2009 Test Collection** by Prof. Florina, Bulik, Chiriac, and Zard, Verónica. Description: CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property. The CLEF-IP track was established in 2009 to investigate IR techniques for patent retrieval and is part of the CLEF 2009 evaluation campaign. The track utilizes a collection of more than 1M patent documents derived from EPO (European Patent Office) sources. The collection contains 4... Updated on November 05, 2021.
- The CLEF-IP 2011 Test Collection** by Prof. Florina, Hertzberg, and Zard, Verónica. Description: CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property. The CLEF-IP track ran from 2009 to 2013 and aimed to investigate IR techniques for patent retrieval. The track utilizes a collection of more than 1.3M patent documents (2.0 million files) derived from EPO (European Patent Office) sources and EuroPCT Applications from four APOs (US, JP, AU, and CA). Updated on November 05, 2021.
- The CLEF-IP 2010 Test Collection** by Prof. Florina, and Zard, Verónica. Description: CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property. The CLEF-IP track was established in 2009 to investigate IR techniques for patent retrieval and is a part of the CLEF 2010 evaluation campaign. The track utilizes a collection of more than 1.3M patent documents (1-2 million files) derived from EPO (European Patent Office) sources, and... Updated on November 05, 2021.
- The CLEF-IP 2013 Test Collection** by Prof. Florina, Hertzberg, and Zard, Verónica. Description: CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property. The CLEF-IP track ran from 2009 to 2013 and aimed to investigate IR techniques for patent retrieval. The track utilizes a collection of more than 1.3M patent documents (2.0 million files) derived from EPO (European Patent Office) sources and EuroPCT Applications from four APOs (US, JP, AU, and CA). Updated on November 05, 2021.
- The CLEF-IP 2012 Test Collection** by Prof. Florina, Hertzberg, and Zard, Verónica. Description: CLEF-IP: Cross-Language Evaluation Forum - Intellectual Property. The CLEF-IP track ran from 2009 to 2013 and aimed to investigate IR techniques for patent retrieval. The track utilizes a collection of more than 1.3M patent documents (2.0 million files) derived from EPO (European Patent Office) sources and EuroPCT Applications from four APOs (US, JP, AU, and CA). Updated on November 05, 2021.
- The MARCIRREC data set** by Prof. Florina. Description: MARCIRREC: The MARCIRREC Dataset - Collection of The Information Retrieval Conference Collection. MARCIRREC is a multi-collection of over 18 million patent applications and granted patents in a unified format, extracted from EP, WO, US, and JP sources, spanning a range from 1978 to June 2008. MARCIRREC is intended as raw material for rese... Updated on November 05, 2021.

At the bottom of the page, there is a footer with logos for 'powered by INVENIO ROOM' and 'enabled by FAIR DATA AUSTRIA', along with links for 'POLICIES', 'TERMS OF USE', 'DATA PROTECTION DECLARATION', and 'CONTACT'.

Search - Workflow



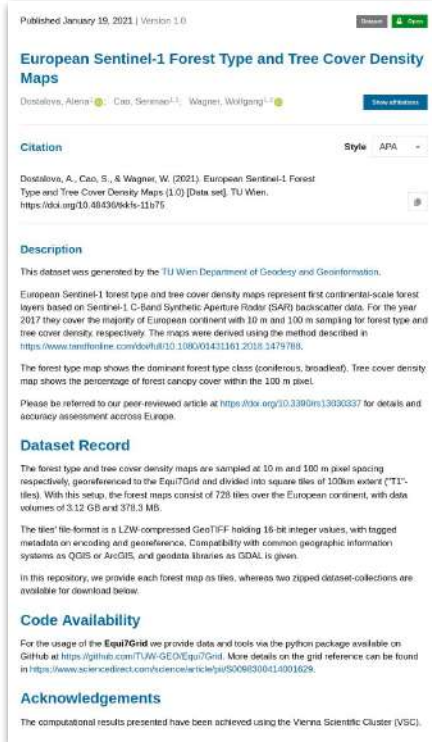
Record Landing Page

The screenshot shows a record landing page for the dataset 'European Sentinel's Forest Type and Tree Cover Density Maps'. The page is structured as follows:

- Header:** Includes 'Home' and 'My Dashboard' links, and the TU logo.
- Dataset Title:** 'European Sentinel's Forest Type and Tree Cover Density Maps'.
- Metadata:** Lists authors (Ferdinand, A., C. A. R., & N. P.), year (2023), and a DOI link.
- Abstract:** A paragraph describing the dataset's purpose and content.
- Dataset Record:** A section detailing the dataset's origin, processing, and availability.
- Code Availability:** A section indicating that the code is available for download.
- Acknowledgements:** A section acknowledging the funding source.
- Files:** A table listing the dataset files with their names and sizes.
- More:** A section with links for 'Download' and 'View'.
- Additional Reads:** A section with links for 'Related records', 'Cite this record', and 'References'.

File Name	Size
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB
EU_Sentinel_FT_Cover_Density_Maps_2023.zip	1.0 GB

Record Landing Page



Published January 19, 2021 | Version 1.0

European Sentinel-1 Forest Type and Tree Cover Density Maps

Dostalova, Alena; Cao, Senqiao; Wagner, Wolfgang

Citation Style: APA

Dostalova, A., Cao, S., & Wagner, W. (2021). European Sentinel-1 Forest Type and Tree Cover Density Maps (1.0) [Data set]. TU Wien. <https://doi.org/10.48430/444-11079>

Description

This dataset was generated by the TU Wien Department of Geodesy and Geoinformation.

European Sentinel-1 forest type and tree cover density maps represent first continental-scale forest layers based on Sentinel-1 C-Band Synthetic Aperture Radar (SAR) backscatter data. For the year 2017 they cover the majority of European continent with 10 m and 100 m sampling for forest type and tree cover density, respectively. The maps were derived using the method described in https://www.semanticscholar.org/paper/10.1007/978-3-319-12038-1_1479788.

The forest type map shows the dominant forest type class (coniferous, broadleaf). Tree cover density map shows the percentage of forest canopy cover within the 100 m pixel.

Please be referred to our peer-reviewed article at <https://doi.org/10.3390/rs13030337> for details and accuracy assessment across Europe.

Dataset Record

The forest type and tree cover density maps are sampled at 10 m and 100 m pixel spacing respectively, georeferenced to the Equi7Grid and divided into square tiles of 100km extent ("T1"-tiles). With this setup, the forest maps consist of 728 files over the European continent, with data volumes of 3.12 GB and 378.3 MB.

The tiles' file format is a LZW-compressed GeoTIFF holding 16-bit integer values, with tagged metadata on encoding and georeference. Compatibility with common geographic information systems as QGIS or ArcGIS, and geodata libraries as GDAL is given.

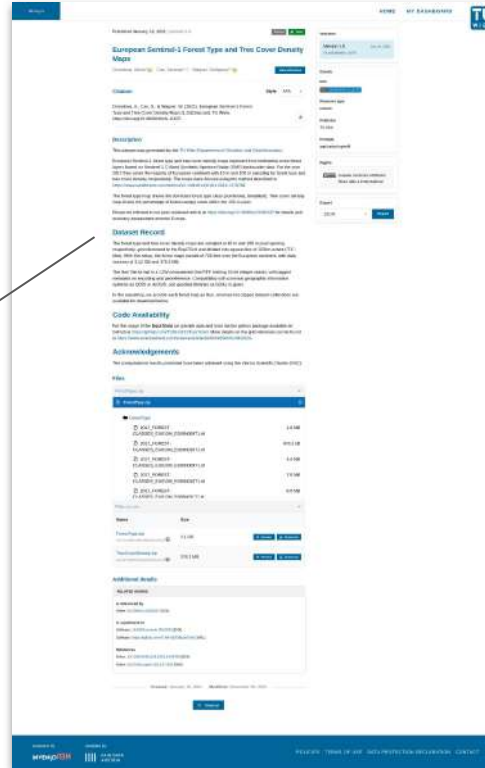
In this repository, we provide each forest map as files, whereas two zipped dataset-collections are available for download below.

Code Availability

For the usage of the Equi7Grid we provide data and tools via the python package available on GitHub at <https://github.com/TUW-GEODEpy7Grid>. More details on the grid reference can be found in <https://www.sciencedirect.com/science/article/pii/S0083300418001629>.

Acknowledgements

The computational results presented have been achieved using the Vienna Scientific Cluster (VSC).



European Sentinel-1 Forest Type and Tree Cover Density Maps

Dataset Record

The forest type and tree cover density maps are sampled at 10 m and 100 m pixel spacing respectively, georeferenced to the Equi7Grid and divided into square tiles of 100km extent ("T1"-tiles). With this setup, the forest maps consist of 728 files over the European continent, with data volumes of 3.12 GB and 378.3 MB.

The tiles' file format is a LZW-compressed GeoTIFF holding 16-bit integer values, with tagged metadata on encoding and georeference. Compatibility with common geographic information systems as QGIS or ArcGIS, and geodata libraries as GDAL is given.

In this repository, we provide each forest map as files, whereas two zipped dataset-collections are available for download below.

Code Availability

For the usage of the Equi7Grid we provide data and tools via the python package available on GitHub at <https://github.com/TUW-GEODEpy7Grid>. More details on the grid reference can be found in <https://www.sciencedirect.com/science/article/pii/S0083300418001629>.

Acknowledgements

The computational results presented have been achieved using the Vienna Scientific Cluster (VSC).

File	Size
10m_FT.tif	1.12 GB
100m_FT.tif	378.3 MB
10m_TCD.tif	1.12 GB
100m_TCD.tif	378.3 MB

Record Landing Page

Published January 19, 2021 | Version 1.0

European Sentinel-1 Forest Type and Tree Cover Density Maps

Dostalova, Alena; Cao, Senqiao; Wagner, Wolfgang

Citation Style: APA
Dostalova, A., Cao, S., & Wagner, W. (2021). European Sentinel-1 Forest Type and Tree Cover Density Maps (1.0) [Data set]. TU Wien. <https://doi.org/10.48430/446-11075>

Description
This dataset was generated by the TU Wien Department of Geodesy and Geoinformation.

European Sentinel-1 forest type and tree cover density maps represent first continental-scale forest layers based on Sentinel-1 C-Band Synthetic Aperture Radar (SAR) backscatter data. For the year 2017 they cover the majority of European continent with 10 m and 100 m sampling for forest type and tree cover density, respectively. The maps were derived using the method described in <https://www.semanticscholar.org/doi/full/10.1108/0004413161120181479788>.

The forest type map shows the dominant forest type class (coniferous, broadleaf). Tree cover density map shows the percentage of forest canopy cover within the 100 m pixel.

Please be referred to our peer-reviewed article at <https://doi.org/10.3390/rs13030337> for details and accuracy assessment across Europe.

Dataset Record
The forest type and tree cover density maps are sampled at 10 m and 100 m pixel spacing respectively, georeferenced to the Equi7Grid and divided into square files of 100km extent ("11"-files). With this setup, the forest maps consist of 728 files over the European continent, with data volumes of 3.12 GB and 378.3 MB.

The files' file-format is a LZW-compressed GeoTIFF holding 16-bit integer values, with tagged metadata on encoding and georeference. Compatibility with common geographic information systems as QGIS or ArcGIS, and geodata libraries as GDAL is given.

In this repository, we provide each forest map as files, whereas two zipped dataset-collections are available for download below.

Code Availability
For the usage of the Equi7Grid we provide data and tools via the python package available on GitHub at <https://github.com/TUW-GEODE/eq7Grid>. More details on the grid reference can be found in <https://www.sciencedirect.com/science/article/pii/S0083304118001629>.

Acknowledgements
The computational results presented have been achieved using the Vienna Scientific Cluster (VSC).

Dataset Record
The forest type and tree cover density maps are sampled at 10 m and 100 m pixel spacing respectively, georeferenced to the Equi7Grid and divided into square files of 100km extent ("11"-files). With this setup, the forest maps consist of 728 files over the European continent, with data volumes of 3.12 GB and 378.3 MB.

Code Availability
For the usage of the Equi7Grid we provide data and tools via the python package available on GitHub at <https://github.com/TUW-GEODE/eq7Grid>. More details on the grid reference can be found in <https://www.sciencedirect.com/science/article/pii/S0083304118001629>.

Acknowledgements
The computational results presented have been achieved using the Vienna Scientific Cluster (VSC).

Files

Name	Size
ForestType.zip	3.1 GB
TreeCoverDensity.zip	378.3 MB

Versions
Version 1.0
10.48430/446-11075
Jan 19, 2021

Details
DOI
[10.48430/446-11075](https://doi.org/10.48430/446-11075)

Resource type
Dataset

Publisher
TU Wien

Formats
application/x-gzotif

Rights
Creative Commons Attribution
Share Alike 4.0 International

Export
JSON

Files

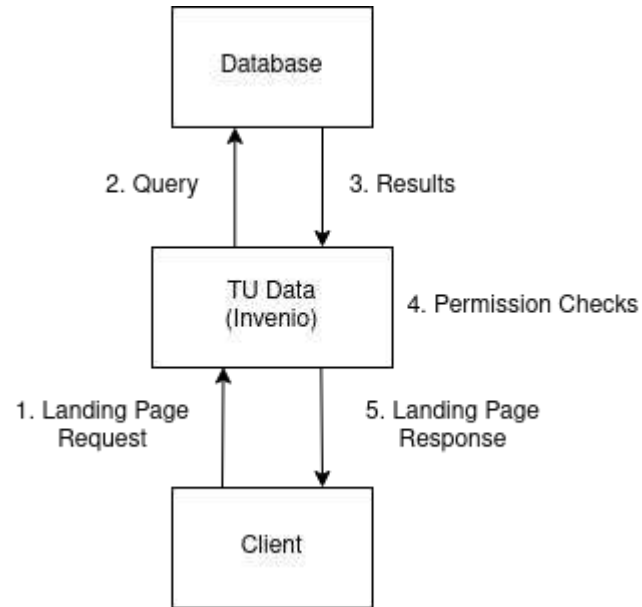
ForestType.zip

- ForestType.zip
- ForestType
- 2017_FOREST-CLASSES_EU10M_E029N009T1.tif (1.0 MB)
- 2017_FOREST-CLASSES_EU10M_E029N107T1.tif (670.2 kB)
- 2017_FOREST-CLASSES_EU10M_E030N009T1.tif (4.4 MB)
- 2017_FOREST-CLASSES_EU10M_E030N107T1.tif (7.5 MB)
- 2017_FOREST-CLASSES_EU10M_E030N111T1.tif (8.0 MB)

Files (15 GB)

Name	Size
ForestType.zip	3.1 GB
TreeCoverDensity.zip	378.3 MB

Record Landing Page - Workflow



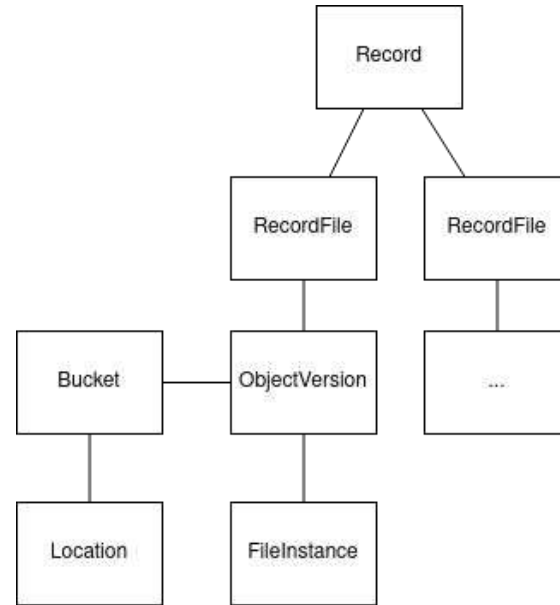
Records: Requirements

- Metadata



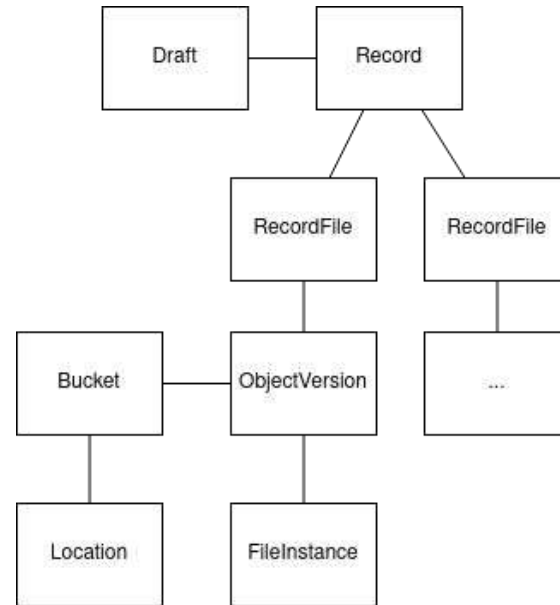
Records: Requirements

- Metadata
- Attached files




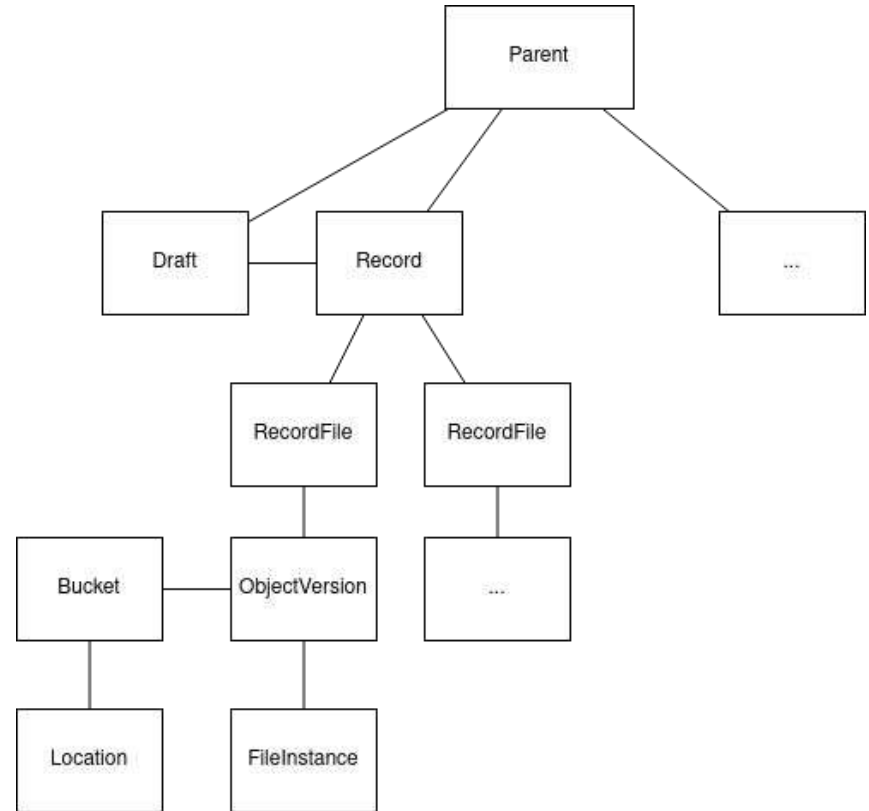
Records: Requirements

- Metadata
- Attached files
- Editable drafts






Records: Requirements

- Metadata
- Attached files
- Editable drafts
- Record versioning
 - PIDs (e.g. DOIs) are only meaningful if the content they point to is fixed
 - What if you found a mistake in your uploaded dataset? 



Records: Metadata

- What does it look like (which standard)?
- How's it stored (format)? 
- Validation 
- Bonus for the Invenio framework: 

It's a library, not an application – keep it flexible!



Image by [DataCite](#)



marshmallow 

Image by [marshmallow-code](#)

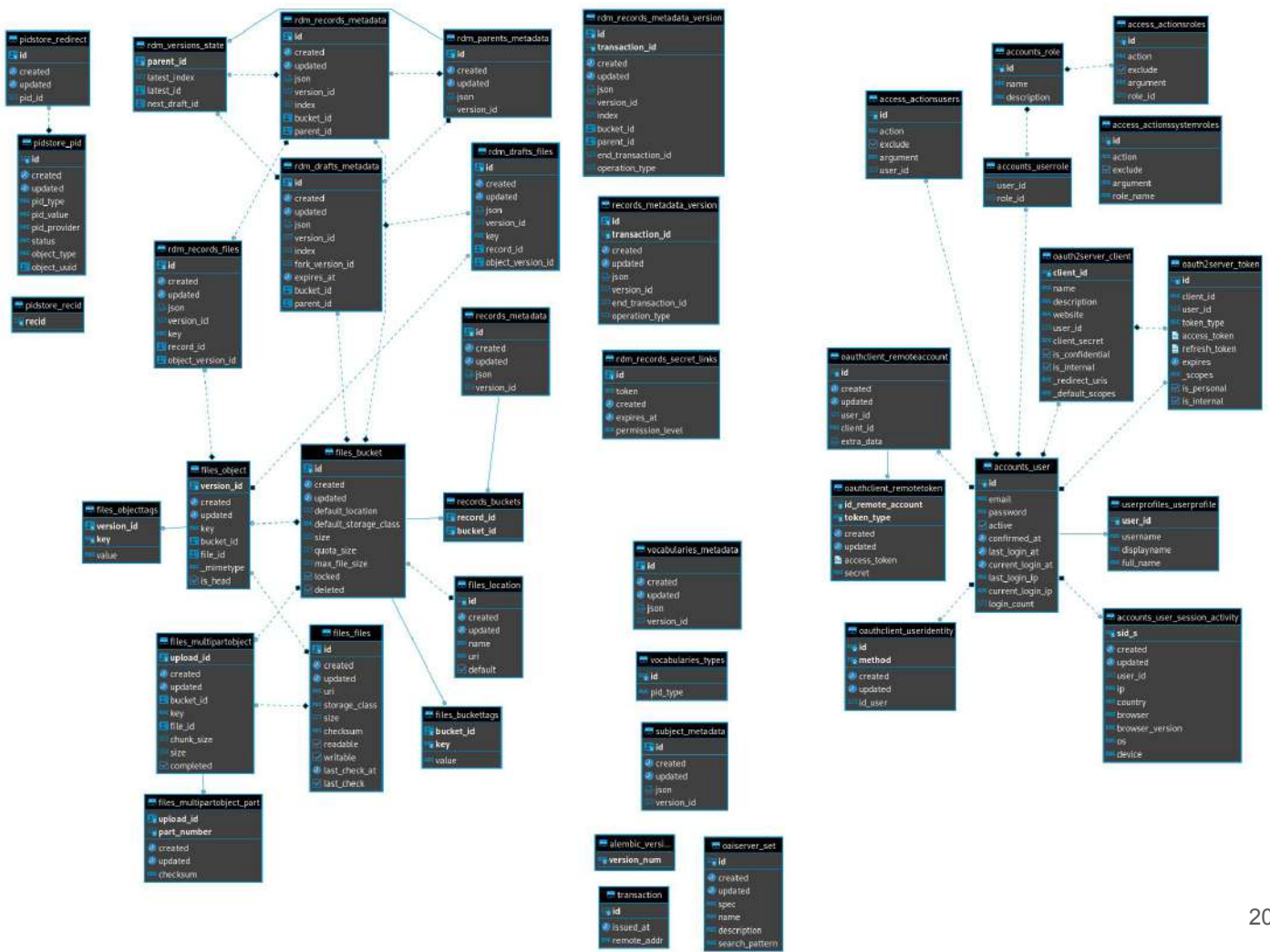
{json}



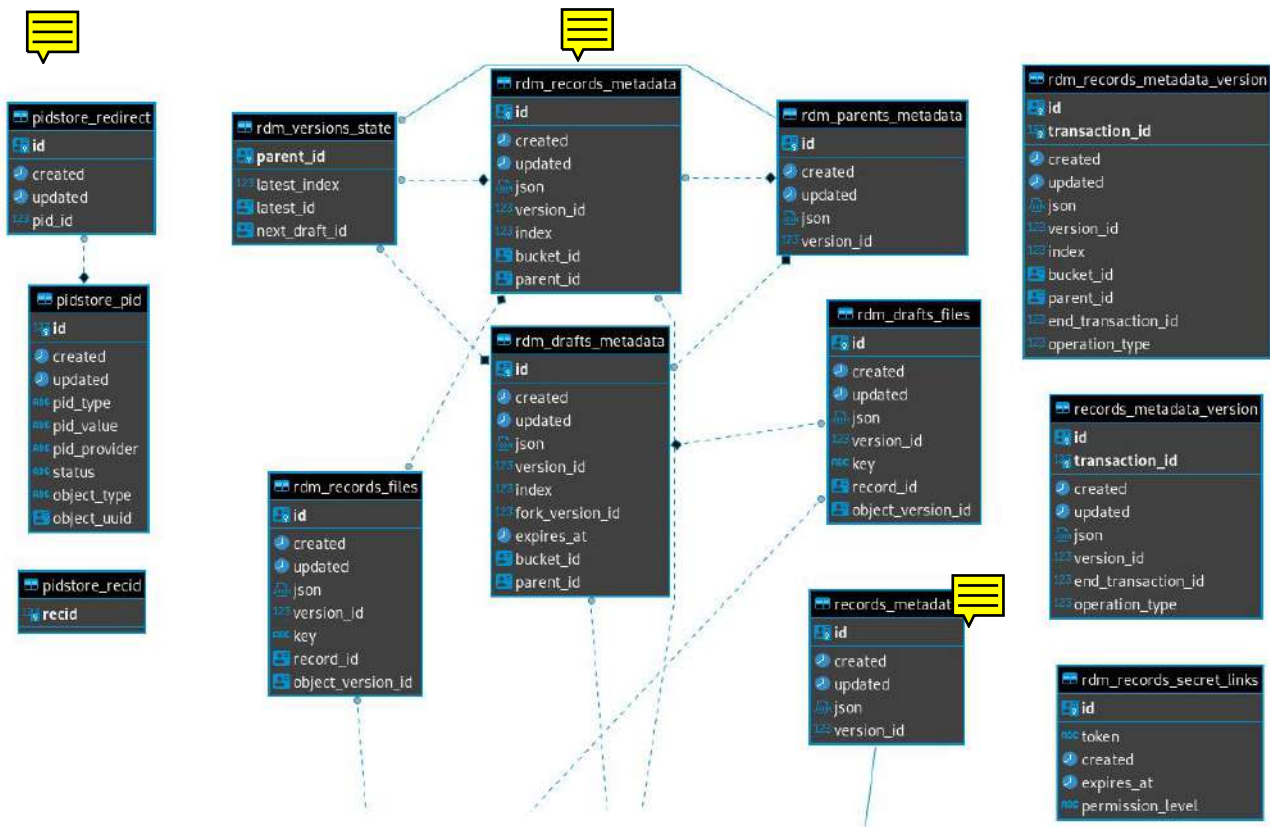
Question: How would you store records?



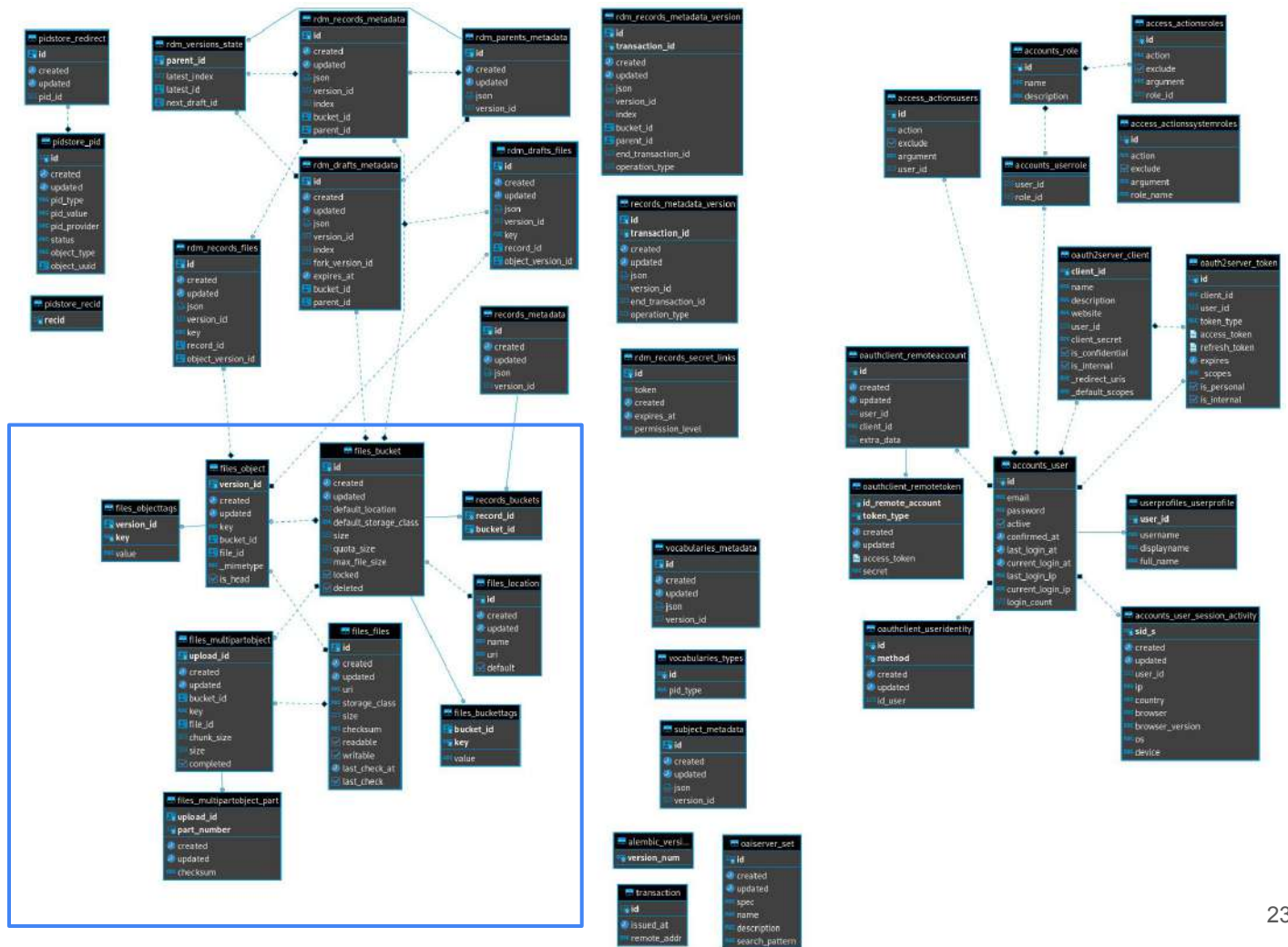
The Database



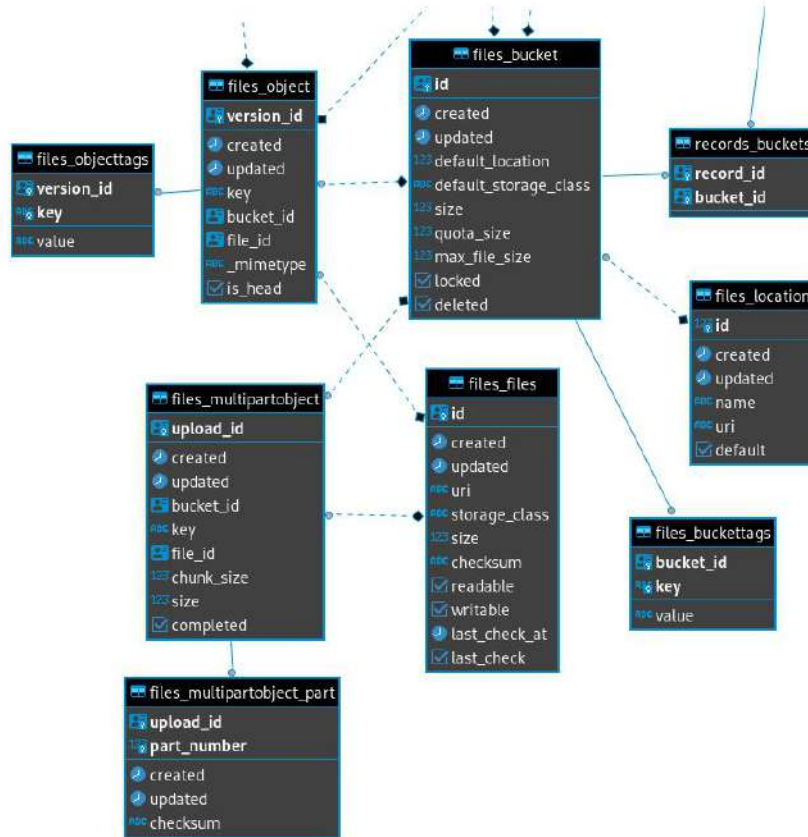
The Database (Records)



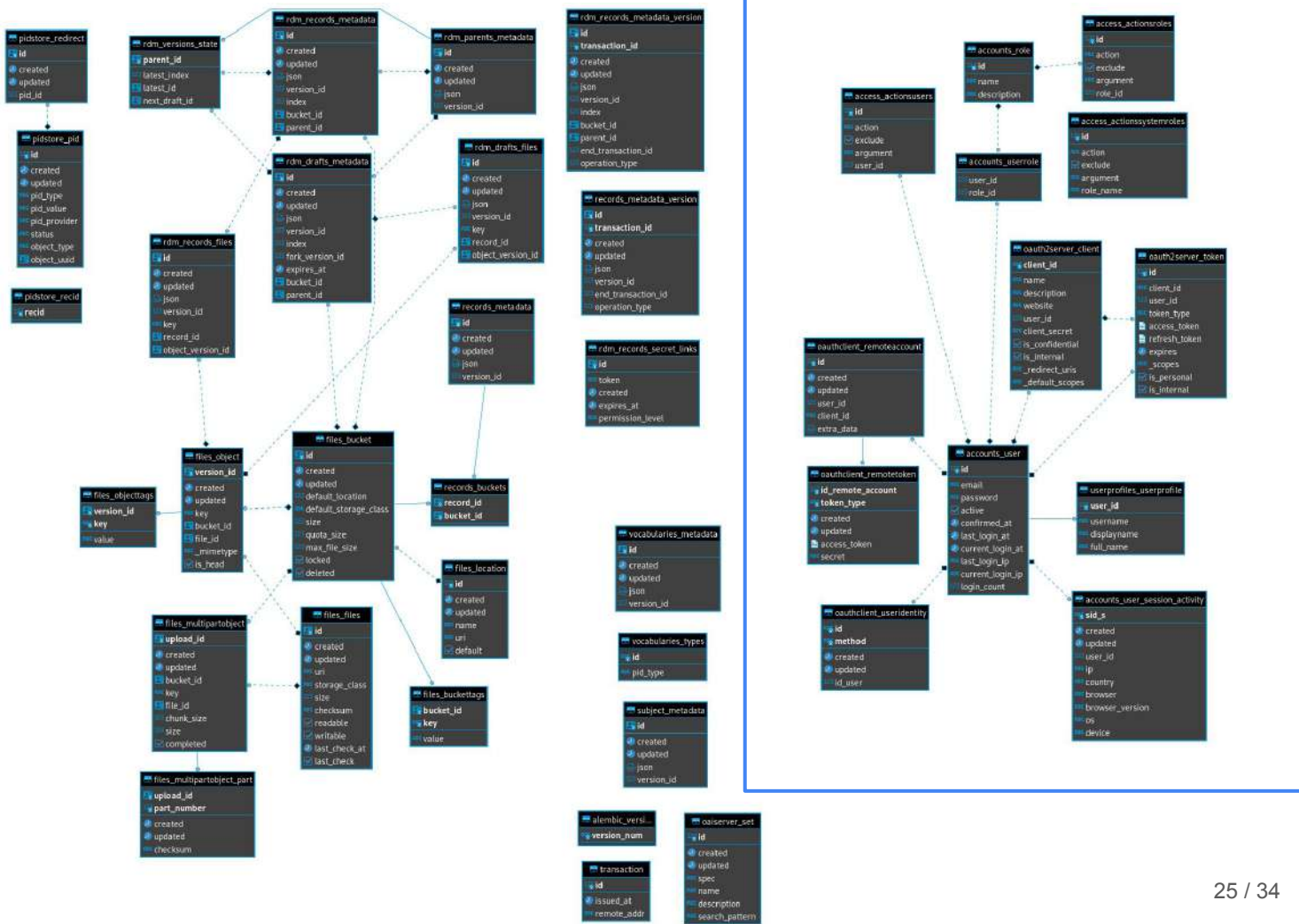
The Database (Files)



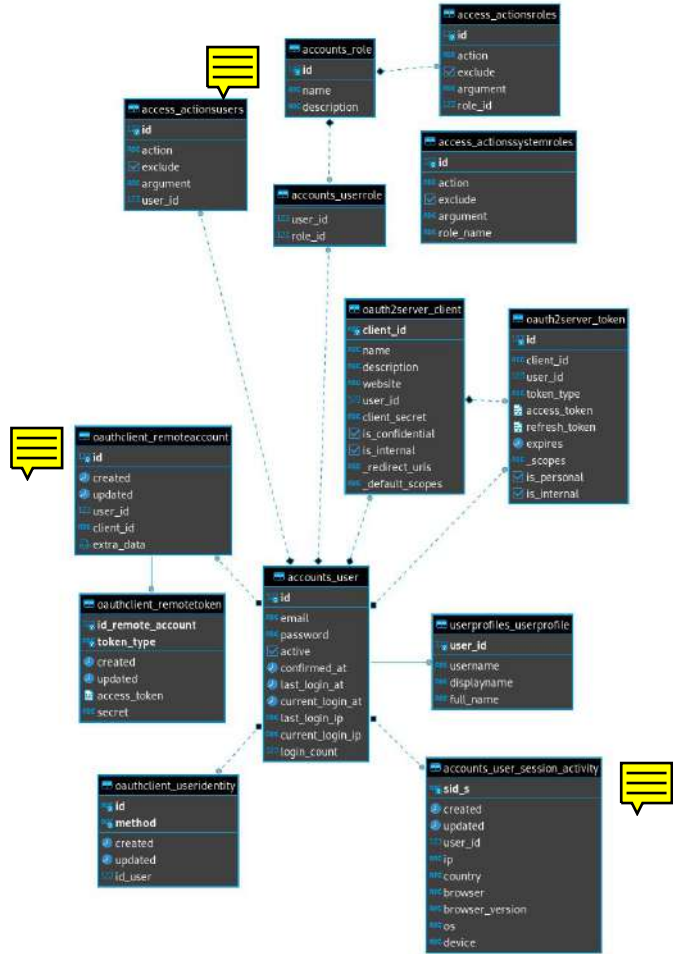
The Database (Files)



The Database (Users)



The Database (Users)



Record Metadata: Why a JSON field instead of ORM?

- A single table with a JSON field
 - Nested structures in JSON
 - No “cluttering” with tables
 - No accidental partial updates



- Fits nicely with ES indexing

- Why not MongoDB then?



- RDBMS still has benefits (ACID, battle-tested, PK lookups are fast, ...)


- Ultimately, just a design decision

rdm_records_metadata	
id	
created	
updated	
json	
version_id	
index	
bucket_id	
parent_id	

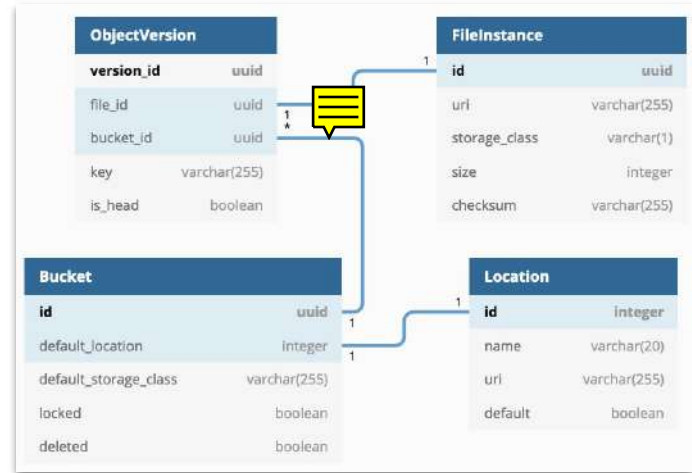
VS .







Records: File Deposits

- Where to store them (local, remote)?
- Make it scalable (like S3: buckets) 
- Fixity checks (checksums)
- Efficiency (e.g. duplicates)
- More information about files in Invenio:



<https://invenio-files-rest.readthedocs.io/en/latest/overview.html>



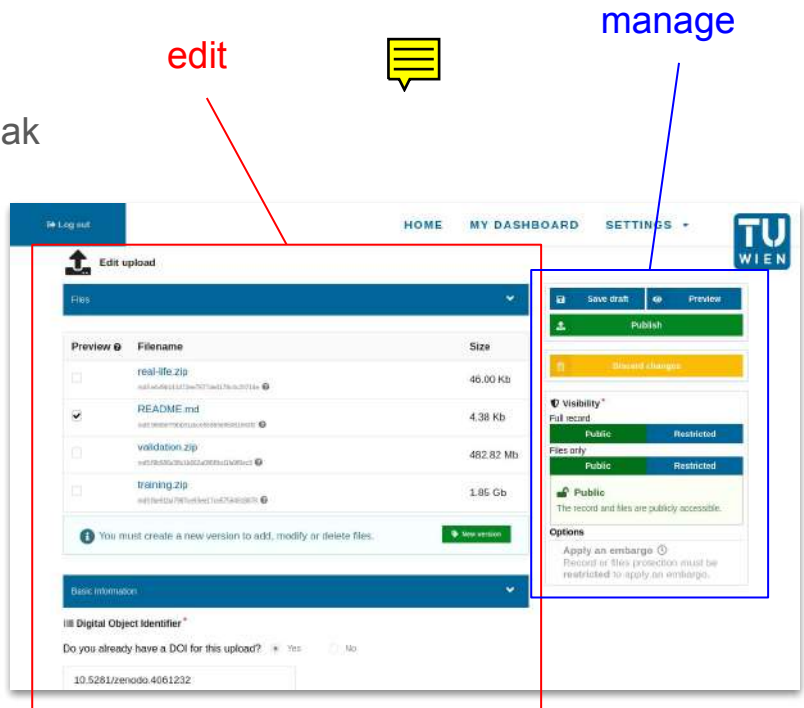
Records: Search

- Map the records to Elasticsearch 
 - May have influence on the [metadata design](#) (e.g. [nested fields](#))
 - Alternatively, add translation layer => increased complexity
- Try to avoid expensive database queries
 - Keep everything relevant in ES 
- Filter out some results (e.g. closed-access records) 
- It's not a database! 
 - Everything can be recreated (re-indexed) from DB at any time

Access Restrictions

- Authentication 
 - Local login vs OAuth2 (resp. OpenID Connect)
 - Get other services integrated, e.g. Keycloak
=> integrates TU SSO
 - Sessions for browsers, tokens for cURL 

- Authorization
 - Different permissions for different actions
 - **read**, **read_files**, **create**, **edit**, **manage**, ...
 - Role-based access (e.g. role “trusted-user”)
 - Share-by-Link (like Google Docs)



The screenshot shows a web interface for file management. At the top right, there are navigation links: "HOME", "MY DASHBOARD", and "SETTINGS". The TU WIEN logo is in the top right corner. The main content area is titled "Edit upload". It features a table of files with columns for "Preview", "Filename", and "Size". The files listed are:

Preview	Filename	Size
<input type="checkbox"/>	real-life.zip	46.00 Kb
<input checked="" type="checkbox"/>	README.md	4.38 Kb
<input type="checkbox"/>	validation.zip	482.82 Mb
<input type="checkbox"/>	training.zip	1.95 Gb

Below the table, there is a message: "You must create a new version to add, modify or delete files." with a "New version" button. At the bottom, there is a section for "Basic information" and a "Digital Object Identifier" field with a dropdown menu and radio buttons for "Yes" and "No". The DOI value is "10.5281/zenodo.4061232".

Annotations on the screenshot:

- A red box highlights the "Edit upload" section, with a red arrow pointing to the word "edit" above it.
- A blue box highlights the "manage" section, with a blue arrow pointing to the word "manage" above it.
- A yellow speech bubble icon is located above the "manage" section.

Share by Link

Can be used to share access to restricted datasets:

The screenshot shows a dataset page with a 'Files' section. A red box highlights a 'Restricted' warning message: "The record is publicly accessible, but files are restricted to users with access." Another red box highlights the 'Files' section, which shows a 'Restricted' status and a message: "The record is publicly accessible, but files are restricted to users with access." The page includes a 'Share' button in the top right corner, which is also highlighted with a red box.

Normal Access

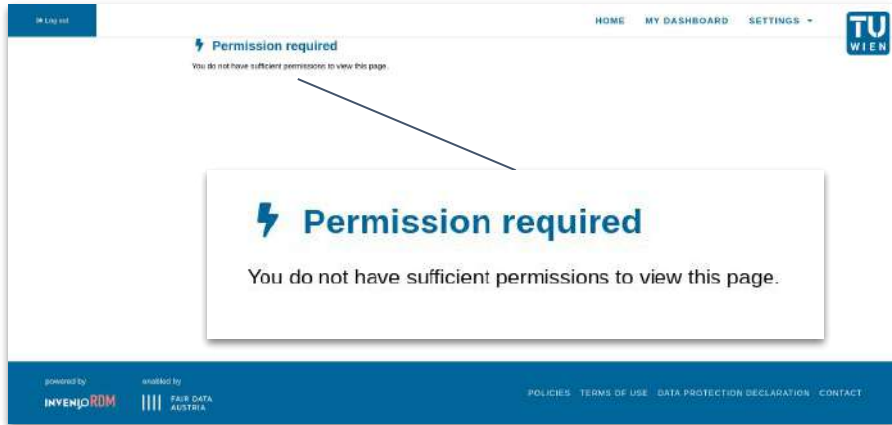
The screenshot shows the same dataset page accessed via a secret link. A red box highlights the 'Share' button in the top right corner. Another red box highlights the 'Files' section, which now shows the full content of the files, including a table of file names and sizes. The table is as follows:

Name	Size	Download	View
HEURISTIC_01	442.0 KB	Download	View
HEURISTIC_02	4.0 KB	Download	View
HEURISTIC_03	442.0 KB	Download	View
HEURISTIC_04	2.0 KB	Download	View

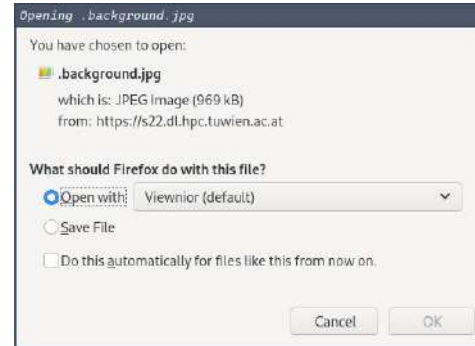
Access via Secret Link

Share by Link

Can also be used to enable double-blind peer reviews:



Metadata cannot be accessed...



... but files can be downloaded

Some Challenges for DevOps

- (Automatic) Deployment
- SSL certificates
- Secrets management
- Migrations (versions, environments)
- Failure detection, recovery, and resistance
- Identifying and fixing bottlenecks and other issues
- Backups
- Licensing
- Life cycles of software
 - Software hits EOL
 - Packages are abandoned
 - New versions introduce breaking changes
- Customizations
 - Custom logic (e.g. for permissions and integration with other services)
 - Custom styling to fit the corporate design
- Etc.

Thanks!

Play with our staging instance:

<https://test.researchdata.tuwien.ac.at>

And send bug reports to:

maximilian.moser@tuwien.ac.at

More documentation:

Invenio RDM: <https://inveniordm.docs.cern.ch/>

Invenio Framework: <https://invenio.readthedocs.io/en/latest/>

Repository systems 2

Tomasz Miksa

Previous lecture



A screenshot of the TU Data website. The header is dark blue with the TU WIEN logo on the left and the text 'TU Data' on the right. Below the header is a search bar with the placeholder text 'Type and press enter to search:'. The main content area is white and features a 'Welcome!' section with a paragraph of text. Below this are four icons in blue boxes: a cloud with an upward arrow, a magnifying glass, a document with stars, and a circular arrow. Each icon has a corresponding text description below it. At the bottom, there is a blue footer with three columns of links: 'TU Data' (Policy, Terms of Use, Privacy Statement, Contact), 'TU Wien' (Center for RDM, RDM Policy), and 'FAIR Data Austria' (FAIR Data Austria).

External visibility

Many institutions only tick off a box 'we have a repository'

'Having a repository' is not enough

- Contents must be discoverable and FAIR
 - Integration with hubs
 - Metadata following standards and machine-actionability



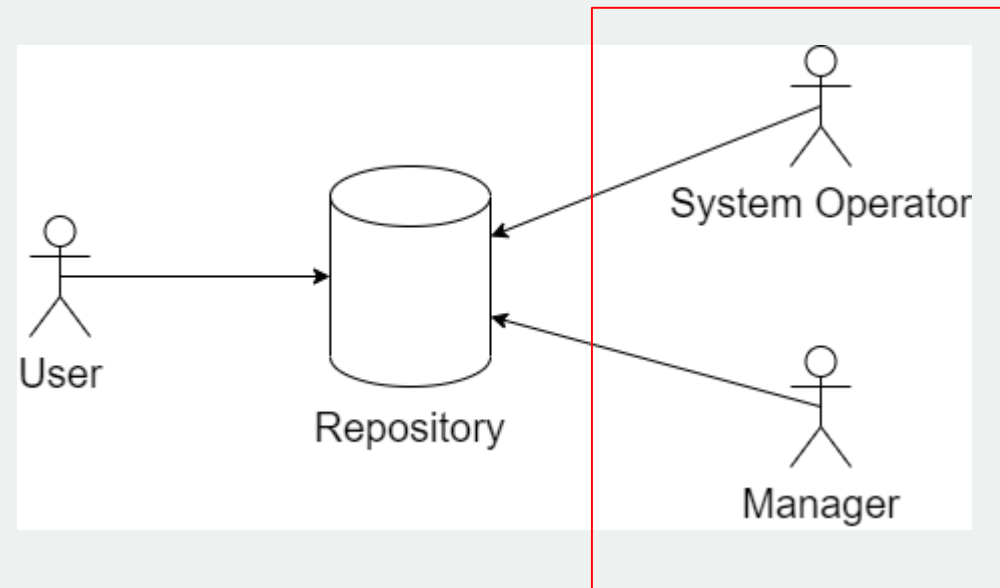
Cargo cult



Agenda

How to make repository contents visible?

- DOI registration bodies
- Interoperability protocols
- Scholix
- DCAT
- Schema.org
- Repository registries



FAIR - Findable

Main focus today:

- **F4. (Meta)data are registered or indexed in a searchable resource**
- Repository itself is a 'searchable resource', but this is not enough


Repositories support also other principles

- e.g. handle access, assign identifiers, etc.
- Not discussed today (and you should know this by now anyway)



DOI registration bodies

DataCite

- DOI registration body
 - Handle based identifier
 - Metadata
 - Provides a range of identifiers to repositories
- Supported by public sector, e.g. DCC, CERN, ANDS
- National desk established at the TU Wien 

Crossref

- DOI registration body
- Supported by private sector, e.g. Elsevier



DOI – repository side

Repository must get an account first to be able to mint DOIs

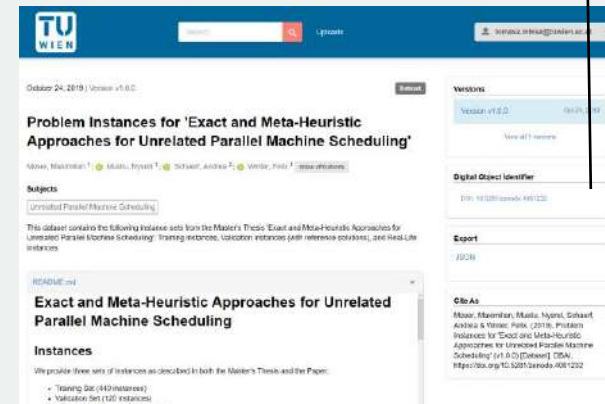
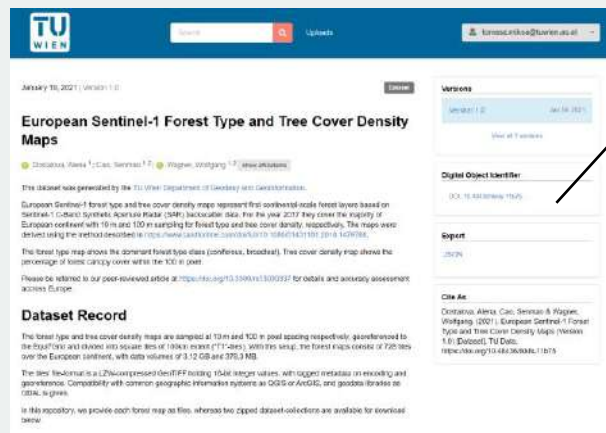
- Prefix: 10.48436
- Suffix: tkkfs-11b75

Digital Object Identifier

DOI: 10.48436/tkkfs-11b75


Digital Object Identifier

DOI: 10.5281/zenodo.4061232



DOI – repository prefix

https://doi.datacite.org/repositories/tuw.tudata/prefixes

DataCite Fabrica 

About Support TUW.TUDATA ▾

TU Data

Info Settings **Prefixes** DOIs

Please ask DataCite Staff if you want to add a prefix.

Type to search... [Search](#)

[Reset All](#)

10.48436

Year created

Year created	Count
<input type="checkbox"/> 2020	1







Added
December 3, 2020, 14:42:17 UTC

About DataCite
What we do
Governance
Members
Steering groups
Staff
Job opportunities

Services
Assign DOIs
Metadata search
Event data
Profiles
re3data
Citation formatter
Statistics
Service status
Contact negotiation

Resources
Metadata schema
Support
Fee Model

Community
Members
Partners
Steering groups
Service providers

Contact us
     
Imprint
Terms and conditions
Privacy policy
● All Systems Operational

[FEEDBACK](#)

DOI – metadata

Repository must provide minimal metadata

- <https://schema.datacite.org/meta/kernel-4.3/>

Table 1: DataCite Mandatory Properties

ID	Property	Obligation
1	Identifier (with mandatory type sub-property)	M
2	Creator (with optional given name, family name, name identifier and affiliation sub-properties)	M
3	Title (with optional type sub-properties)	M
4	Publisher	M
5	PublicationYear	M
10	ResourceType (with mandatory general type description sub-property)	M

Table 2: DataCite Recommended and Optional Properties

ID	Property	Obligation
6	Subject (with scheme sub-property)	R
7	Contributor (with optional given name, family name, name identifier and affiliation sub-properties)	R
8	Date (with type sub-property)	R
9	Language	O
11	AlternateIdentifier (with type sub-property)	O
12	RelatedIdentifier (with type and relation type sub-properties)	R
13	Size	O
14	Format	O
15	Version	O
16	Rights	O
17	Description (with type sub-property)	R
18	GeoLocation (with point, box and polygon sub-properties)	R
19	FundingReference (with name, identifier, and award related sub-properties)	O

TU Data

[Info](#) [Settings](#) [Prefixes](#) [DOIs](#)

Create DOI (Form)

More information about DOI registration via form can be found on our [Support Website](#). Required properties are marked with a red asterix.

Required Properties

* **DOI** A globally unique string that identifies the resource and can't be changed.

10.48436 z0fz-p653

Click the circle icon for a new random suffix, or the cross icon to delete the random suffix and enter a value manually.

* **State** The state determines whether a DOI is registered and findable. Once in Registered or Findable state, a DOI can't be set back to Draft state. [More ...](#)

- Draft only visible in Fabrica, DOI can be deleted
- Registered registered with the DOI Resolver
- Findable registered with the DOI Resolver and **indexed in DataCite Search**

* **URL** The location of the landing page with more information about the resource.

URL

Should be a https URL – within the allowed domain(s) of your repository if domain restrictions are enabled in the repository settings. Http and ftp are also supported.

* **Creators** The main researchers or organizations involved in producing the resource, in priority order.

Name Identifier

Name Identifier

Uniquely identifies an individual or legal entity, according to various schemas, e.g. ORCID, ROR or ISNI. Use name identifier expressed as URL. The Given Name, Family Name and Name will automatically be filled out for ORCID and ROR identifiers.

+ Add another name identifier

- Person Organization Unknown

+ Add another creator Hide 1 creator

*** Titles** One or more names or titles by which the resource is known.

Title

Title Type

Select Title Type

Language

Select Language

+ Add another title Hide 1 title

*** Publisher** The name of the entity that holds, archives, publishes prints, distributes, releases, issues, or produces the resource.

Publisher

This property will be used to formulate the citation, so consider the prominence of the role.

*** Publication Year** The year when the resource was or will be made publicly available.

Publication Year

Must be a year between 1000 and 2021.

*** Resource Type General** The general type of the resource.

Select Resource Type General

- Audiovisual
- Book
- Book chapter
- Collection
- Computational notebook
- Conference paper
- Conference proceeding
- Data paper
- Dataset
- Dissertation

Recommended Properties

Subjects Subject, keyword, classification code, or key phrase describing the resource.

+ Add subject

Contributors The institution or person responsible for collecting, managing, distributing, or otherwise contributing to the development of the resource.

DOI – suffix

<https://help.zenodo.org>

- Why do you include “zenodo” in the DOI?
- Why don't the DOIs have a version number suffix like “.v1”?

Cool DOIs

- <https://blog.datacite.org/cool-dois/>

Digital Object Identifier

DOI: 10.48436/tkkfs-11b75

Digital Object Identifier

DOI: 10.5281/zenodo.4061232

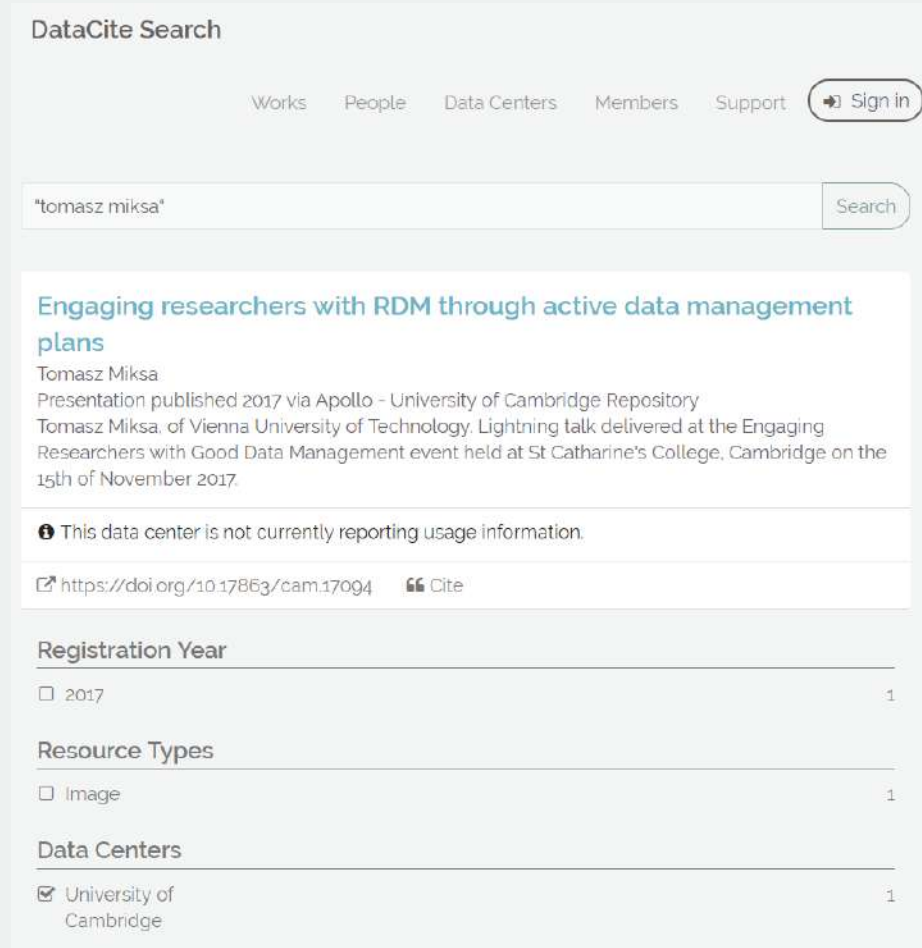


DOI registration bodies - example

Example shows

- Document located at Cambridge Repository
- Document has a DOI
- DOI was minted by Data Cite
- Data Cite has metadata about each DOI

Data Cite provides access to metadata registry



The screenshot shows the DataCite Search interface. At the top, there are navigation links for Works, People, Data Centers, Members, and Support, along with a Sign in button. A search bar contains the text "tomasz miksa" and a Search button. Below the search bar, a search result is displayed for the document "Engaging researchers with RDM through active data management plans" by Tomasz Miksa. The result includes a description of the presentation published in 2017 at the University of Cambridge Repository. Below the description, there is a note indicating that the data center is not currently reporting usage information. The DOI link is provided as <https://doi.org/10.17863/cam.17094>. Below the DOI link, there are three filter sections: "Registration Year" with a checkbox for 2017 (count 1), "Resource Types" with a checkbox for Image (count 1), and "Data Centers" with a checked checkbox for University of Cambridge (count 1).


<https://search.datacite.org>

Interoperability protocols

OAI-PMH - Open Archives Initiative Protocol for Metadata Harvesting

- query to discover repository contents
- only for metadata
- not for depositing 

SWORD - Simple Web-service Offering Repository Deposit

- deposit to multiple repositories at once 
- deposit by third party systems (e.g. lab equipment)

OAI-PMH

Data remains within a repository

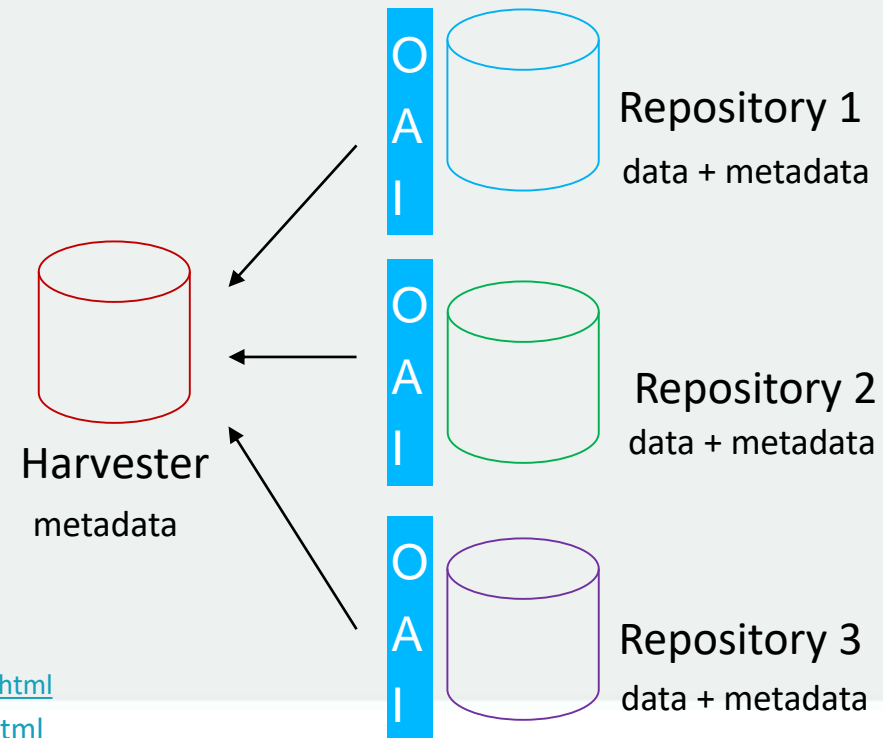
Harvester aggregates metadata

Useful to aggregate data

- e.g. for domain, country, institution

Dublin Core by default

- Other Metadata standards can be added
 - OpenAIRE requires DataCite



https://guidelines.openaire.eu/en/latest/data/use_of_oai_pmh.html

<http://www.openarchives.org/OAI/openarchivesprotocol.html>

OAI-PMH in OpenDOAR (registry)


The screenshot shows the OpenDOAR interface for the University of Vienna PHAIDRA repository. The page includes a navigation menu with links for About, Search, Statistics, Policy Support, Our Work, Contact, and Admin. The repository information is displayed in a table format.

Repository Information	
Repository Name	University of Vienna PHAIDRA <small>(English)</small>
Repository Type	Institutional
Description	This site provides access to the digitised copies of the institution's collection as well as teaching material. The interface is available in German, English, Italian and Serbian.
Repository URL	https://phaidra.univie.ac.at
OAI-PMH URL	https://services.phaidra.univie.ac.at/api/oai
Year Established	2008
Software Name	Fedora <small>(version 3.2)</small>
Languages	English German Italian Serbian
Content Types	Books, Chapters and Sections Learning Objects Other Special Item Types
Subjects	Multidisciplinary
Additional Information	All objects contain a permanent digital signature, and the objects can be described in multiple languages.
Record Count	Metadata: 54229

<https://v2.sherpa.ac.uk/id/repository/1726?template=opendoar>

OAI-PMH

validator.oaipmh.com/#Identify



Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH)
Validator & data extractor Tool

Download and evaluate XML metadata from **OAI-PMH enabled** digital libraries.

Validate URL ✓ Validate By Direct Input Download XML </> REST API About References

OAI-PMH URL: Check now »
Example OAI-PMH URL: <https://oai.datacite.org/oai>

AVAILABLE COMMANDS

- Identify
- ListMetadataFormats
- ListSets
- ListIdentifiers
- ListRecords OAI_DC
- ListRecords OAI_OPENAIRE

✓ Identify Validation 7 XML Result 1 KB

- ✓ HTTP status 200
- ⚠ Content type text/xml; charset=utf-8
- ✓ Content XML checked.
- ⚠ Request time too much 5.009
- ✓ XML complies with OAI-PMH XML Schema <http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd>
- ✓ OAI-PMH protocol version is 2.0.
- ✓ Valid adminEmail admin.phaidra@univie.ac.at

<https://validator.oaipmh.com/>

OAI-PMH

OAI-PMH URL: [Check now »](#)

Example OAI-PMH URL: <https://oai.datacite.org/oai>

AVAILABLE COMMANDS

- Identify
- ListMetadataFormats
- ListSets**
- ListIdentifiers
- ListRecords OAI_DC
- ListRecords OAI_OPENAIRE

ListSets help

✓ ListSets Validation **48** XML Result **40 KB**

- ✓ HTTP status 200
- ⚠ Content type text/xml; charset=utf-8
- ✓ Content XML checked.
- ✓ Request time is 0.43 sec
- ✓ XML complies with OAI-PMH XML Schema <http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd>
- ✓ Found set "Nationale und Studienkataloge der Juridischen Fakultät" with setSpec "o688263".
- ✓ Found set "Nachlass von Erwin Schrödinger (1887-1961)" with setSpec "europeana_a1051".
- ✓ Found set "Nachlass Maximilian Hell: Aufzeichnungen des berühmten Wiener Astronomen" with setSpec "o314596".
- ✓ Found set "iPRES 2019 - 16th International Conference on Preservation of Digital Objects: iPRES 2019" with setSpec "ipres2019".
- ✓ Found set "iPRES 2017 - Proceedings of the 14th Conference on Preservation of Digital Objects" with setSpec "ipres2017".
- ✓ Found set "iPRES 2012 - Proceedings of the 9th International Conference on Preservation of Digital Objects: iPRES 2012 - Digital Curation Institute, iSchool, Toronto" with setSpec "ipres2012".
- ✓ Found set "Nachlass von Hans Thirring (1888-1976)" with setSpec "europeana_a1052".
- ✓ Found set "DiFaB, Digitales ForschungsArchiv Byzanz" with setSpec "difab".
- ✓ Found set "Mündliche Geschichtsüberlieferung: Österreichische Exilpublizisten im Widerstand gegen den Nationalsozialismus – Ein Oral Video History-Projekt 1988-1996" with setSpec "o868526".
- ✓ Found set "Archiv der Universität Wien, digitale Objekte" with setSpec "archivunivie".
- ✓ Found set "iPRES 2005 - 2nd International Conference on Preservation of Digital Objects: iPRES 2005 - Göttingen" with setSpec "ipres2005".



OpenAIRE

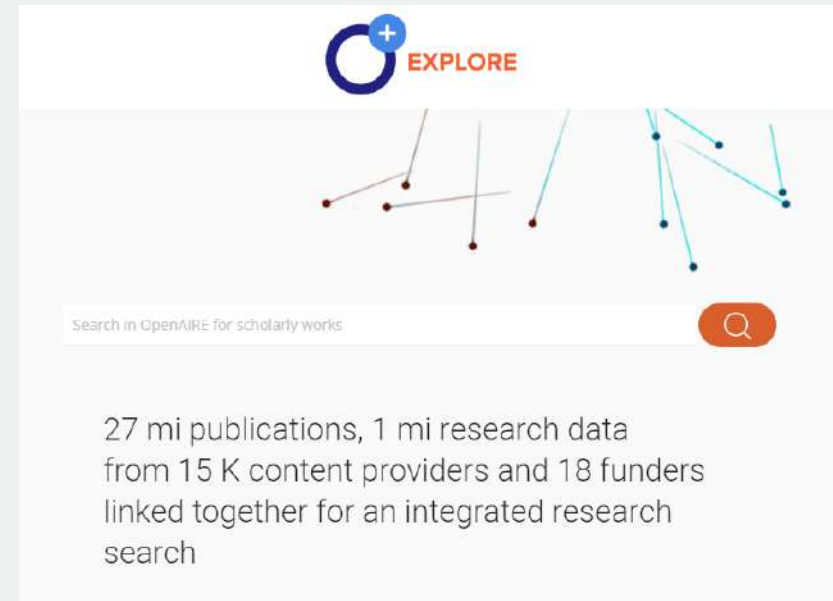
Open Access Infrastructure for Research in Europe

- Launched 2009 by European Commission
- Promotes Open Access

Network of experts

Technical infrastructure

- Harvest research output 
 - Data, publications
 - Link it
 - Monitor
- Zenodo 



<https://explore.openaire.eu>

OpenAIRE

Check what kind of sources are indexed and how metadata is collected

- <https://www.openaire.eu/aggregation-and-content-provision-workflows>

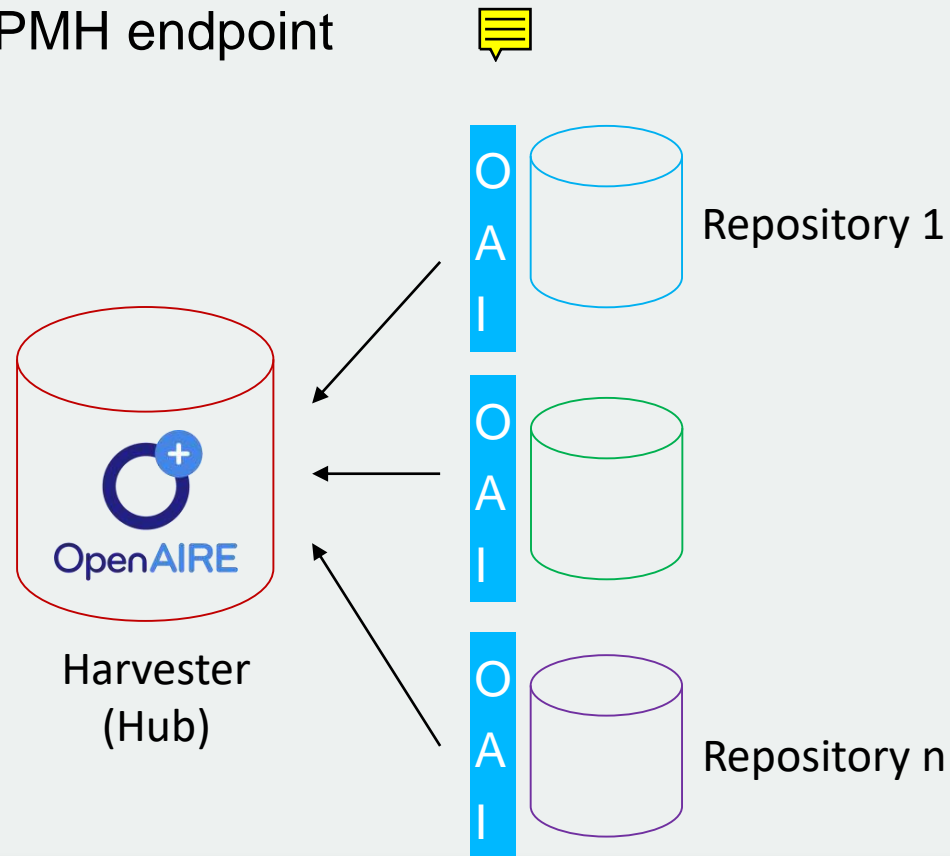
What else openAIRE does

- <https://www.openaire.eu/faqs>

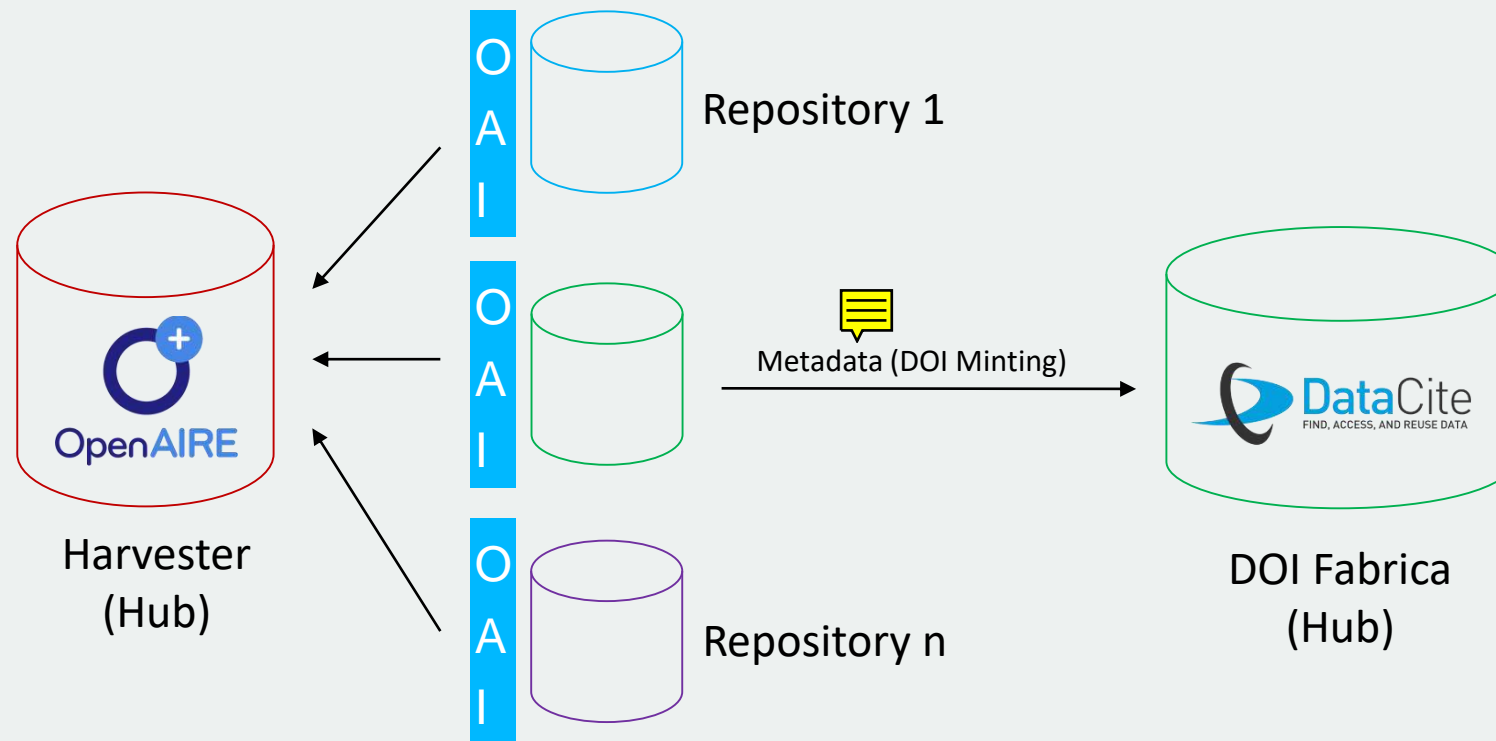


OpenAIRE - compliance

Repository must have an OAI-PMH endpoint
Metadata in DataCite format



OpenAIRE AND/OR DataCite?



OpenAIRE AND/OR DataCite?

Not every repository

- Must participate in OpenAIRE
 - lack of OAI-PMH
 - not located in Europe
- Must mint DOIs
 - Handles, ARKs, etc.



Everyone must decide

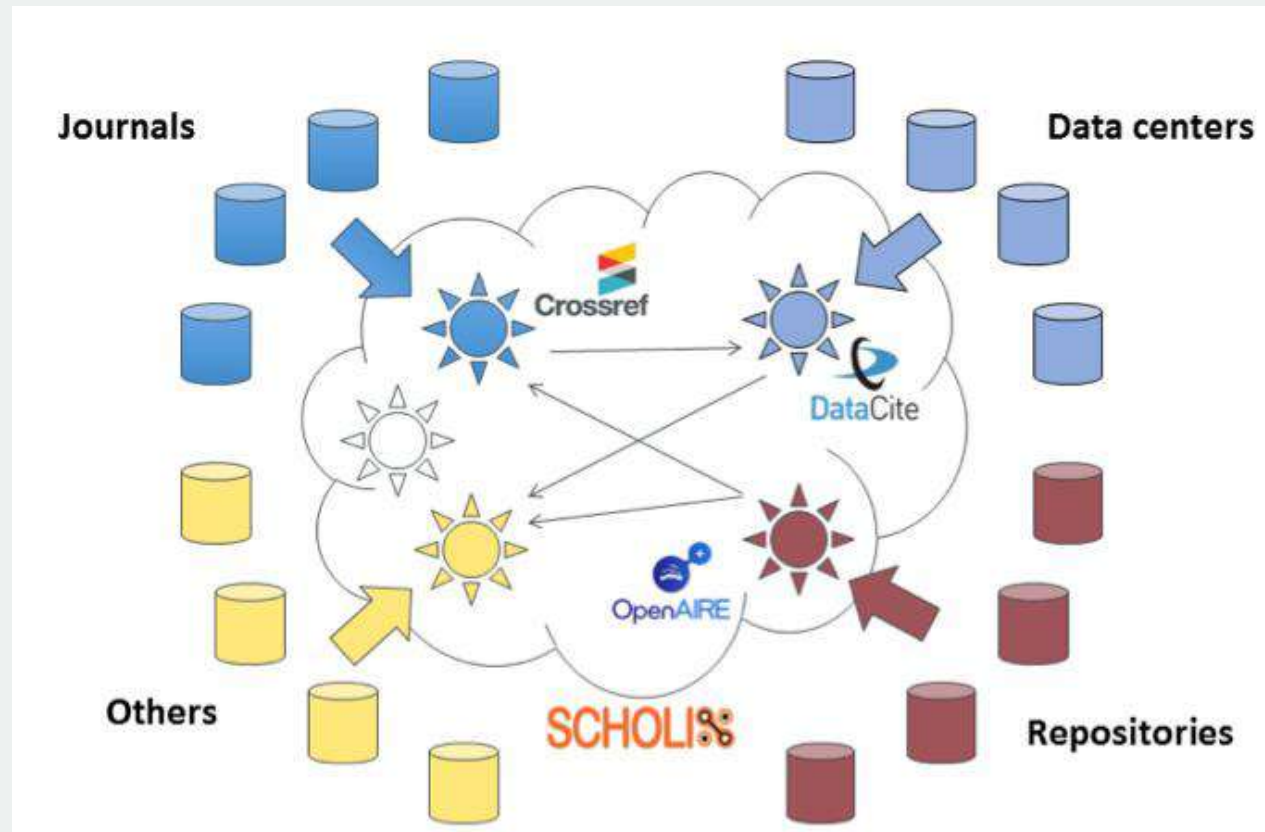
Eventually it theoretically makes little difference...

Scholix



Goal

- exchange links between publications and data



<http://www.dlib.org/dlib/january17/burton/01burton.html>

Scholix

At the core of the conceptual model is the *link* between two *objects*.

Main focus: literature and data.

- theoretically also: software, algorithms, models, protocols, tweets, comments, and so on.
- practically: not for the time being the focus

To become a contributor

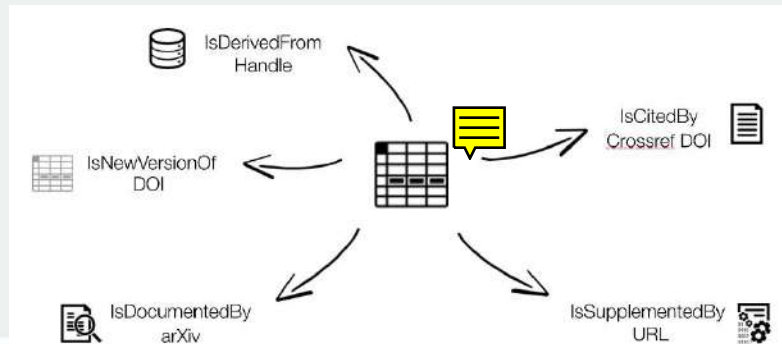
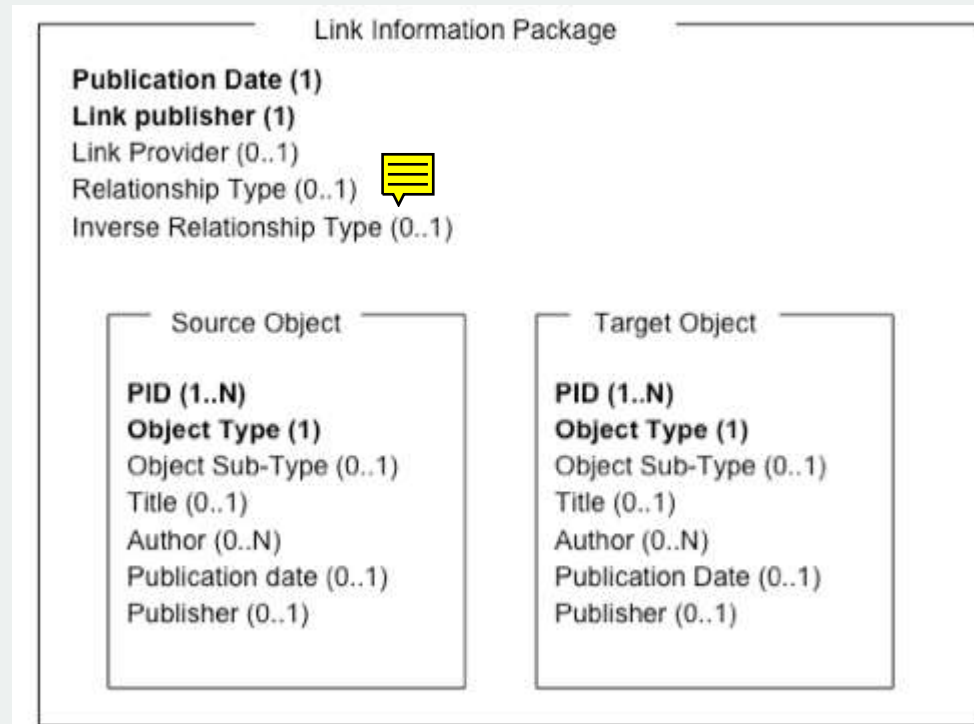
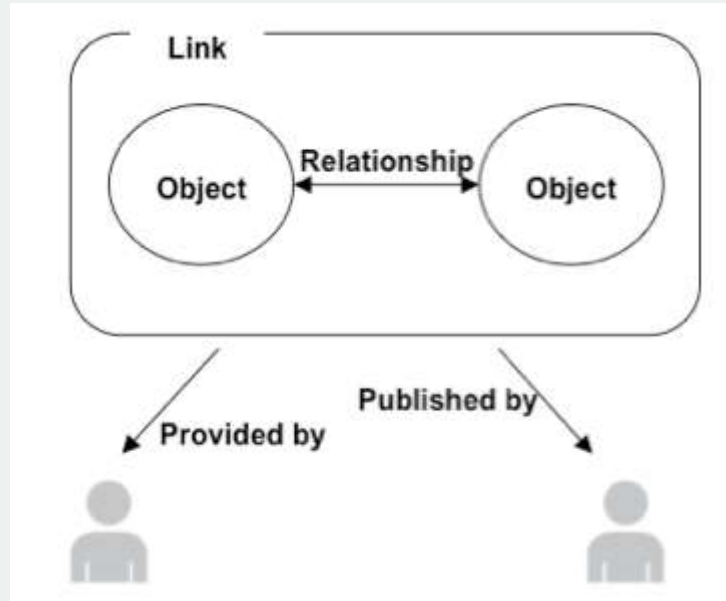
- Feed your data-literature link information to an existing Scholix hub using your existing community standards
 - e.g. OpenAIRE or DataCite registries

To retrieve links

- <http://api.scholexplorer.openaire.eu/v2/ui/>

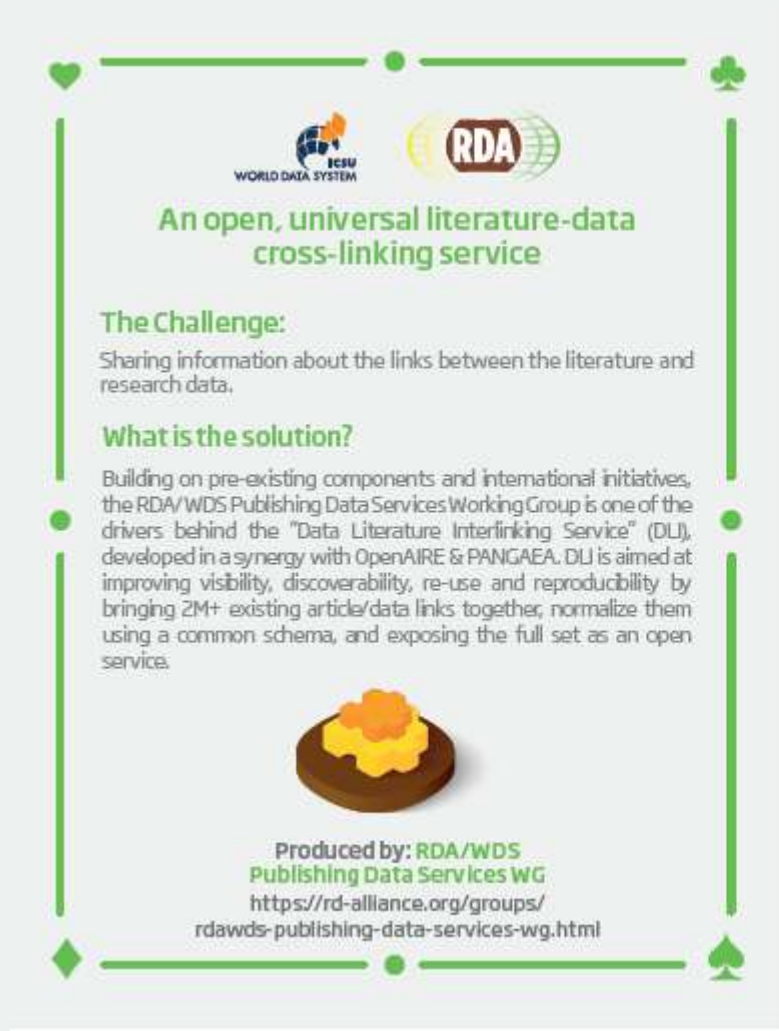


Scholix





<http://api.schoexplorer.openaire.eu/v2/ui/>

Scholix was developed at RDA




The infographic is enclosed in a green border with decorative corner icons: a heart at the top-left, a clover at the top-right, a diamond at the bottom-left, and a spade at the bottom-right. At the top, it features the logos for the ICIV World Data System and RDA. The main title is 'An open, universal literature-data cross-linking service'. Below this, it addresses 'The Challenge' of sharing literature and research data links, and 'What is the solution?' which involves building on existing initiatives like OpenAIRE and PANGAEA to create a 'Data Literature Interlinking Service' (DLI) that normalizes and exposes a large set of article/data links as an open service. At the bottom, it includes a small image of a stack of yellow and orange flowers and provides production credits to the RDA/WDS Publishing Data Services WG along with a URL.

An open, universal literature-data cross-linking service

The Challenge:
Sharing information about the links between the literature and research data.

What is the solution?
Building on pre-existing components and international initiatives, the RDA/WDS Publishing Data Services Working Group is one of the drivers behind the "Data Literature Interlinking Service" (DLI), developed in a synergy with OpenAIRE & PANGAEA. DLI is aimed at improving visibility, discoverability, re-use and reproducibility by bringing 2M+ existing article/data links together, normalize them using a common schema, and exposing the full set as an open service.



Produced by: RDA/WDS
Publishing Data Services WG
<https://rd-alliance.org/groups/rdawds-publishing-data-services-wg.html>



This infographic has a green background and is surrounded by a repeating pattern of the RDA logo. It features two main sections: 'What is the impact?' and 'Find out more about the RDA/WDS Publishing Data Services WG Recommendation'. The impact section describes how large-scale access to literature-data links improves data publishing services and discoverability, and how data centres can better track and present data usage. A QR code is located in the bottom right corner.

What is the impact?

Accessing and using literature-data links at large scale in an efficient and reliable way allows different stakeholders in the data publishing landscape to improve their services, increase data discoverability and usability.

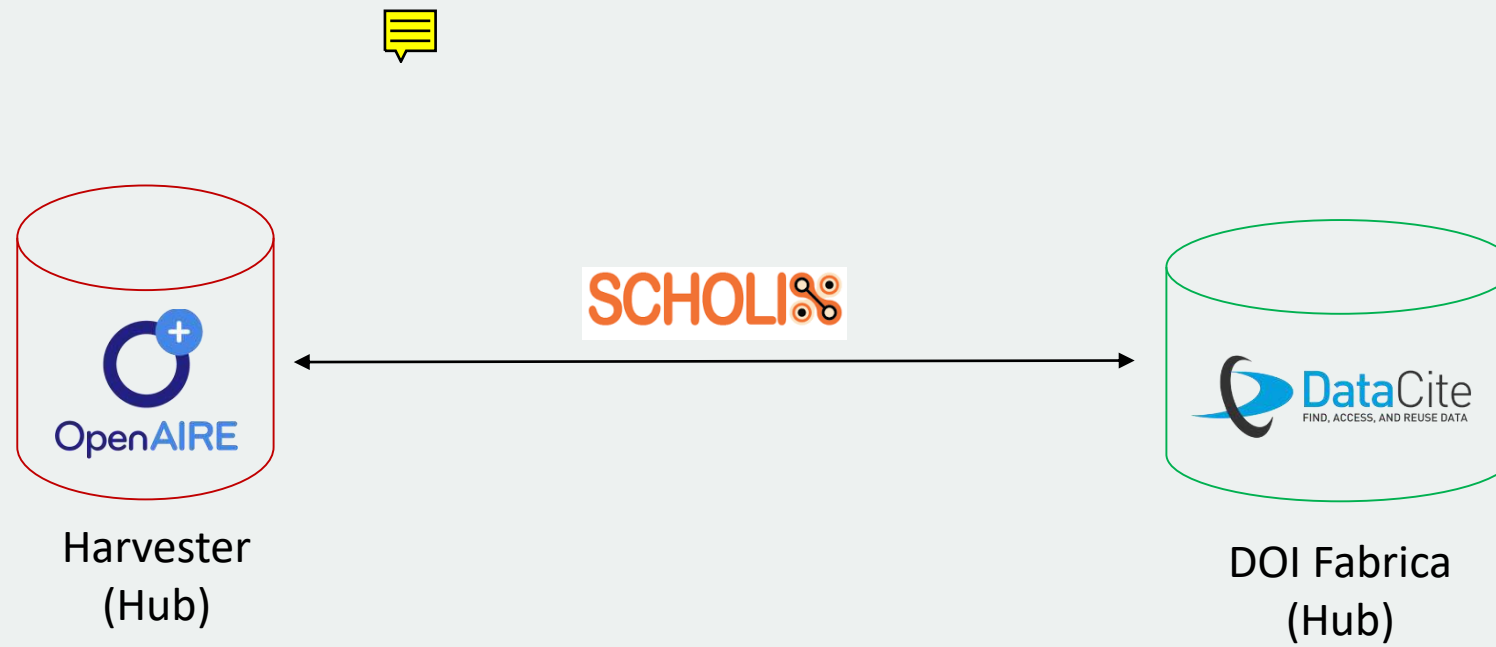
Data centres will be able to assess much better how often their data is used in the literature, and present their users with links to relevant publications.



Find out more about the RDA/WDS
Publishing Data Services WG
Recommendation

<https://www.rd-alliance.org/groups/rdawds-scholarly-link-exchange-scholix-wg>

Scholix



DCAT

W3C Data Catalog Vocabulary (DCAT)

- RDF vocabulary for interoperability between data catalogues
- decentralized publishing
- facilitates federated dataset search

Relaxed constraints

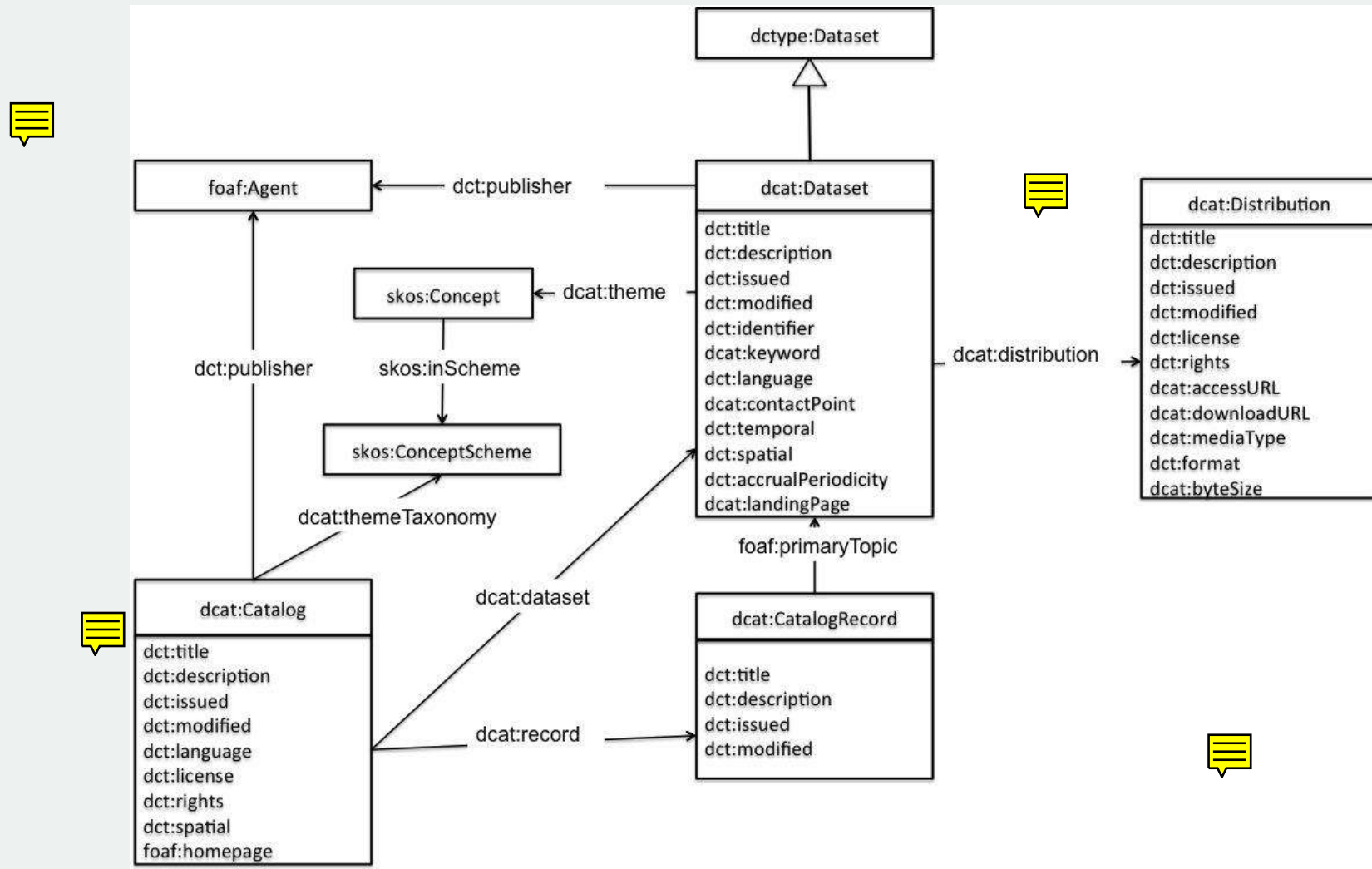
- Most fields are optional

No specific deployment method

- RDF via SPARQL, embedded in HTML, serialised to RDF/XML or Turtle, etc.

Mostly used in open governmental data repositories

DCAT



Specification and examples: <https://www.w3.org/TR/vocab-dcat/>

DCAT Application Profile



DCAT profile is a specification that adds additional constraints
Application Profile for open data in Europe

- <https://joinup.ec.europa.eu/solution/dcat-application-profile-data-portals-europe/about>

4.3. Dataset

4.3.1. Mandatory properties for Dataset

Property	URI	Range	Usage note	Card
description	dct:description	rdfs:Literal	This property contains a free-text account of the Dataset. This property can be repeated for parallel language versions of the description.	1..n
title	dct:title	rdfs:Literal	This property contains a name given to the Dataset. This property can be repeated for parallel language versions of the name.	1..n

4.3.2. Recommended properties for Dataset

Property	URI	Range	Usage note	Card
contact point	dcat:contactPoint	vcard:Kind	This property contains contact information that can be used for sending comments about the Dataset.	0..n
dataset distribution	dcat:distribution	dcat:Distribution	This property links the Dataset to an available Distribution.	0..n
keyword/tag	dcat:keyword	rdfs:Literal	This property contains a keyword or tag describing the Dataset.	0..n
publisher	dct:publisher	foaf:Agent	This property refers to an entity (organisation) responsible for making the Dataset available.	0..1
theme/category	dcat:theme, subproperty dct:subject	of skos:Concept	This property refers to a category of the Dataset. A Dataset may be associated with multiple themes.	0..n

DCAT - example



```
<dcat:Dataset rdf:about="https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe">
  < dct:title>Radabstellplätze Wien</dct:title>
  < dct:description>Radabstellplätze in Wien</dct:description>
  < dcat:keyword>fahrrad</dcat:keyword>
  < dcat:keyword>räder</dcat:keyword>
  < dcat:keyword>verkehr</dcat:keyword>
  < dcat:keyword>fahrräder</dcat:keyword>
  < dcat:distribution>
    < dcat:Distribution rdf:about="https://www.data.gv.at/dataset/af8e02b6-1e03...">
      < dct:title>Radabstellplätze 2016</dct:title>
      < dct:format>CSV</dct:format>
      < dcat:accessURL rdf:resource="https://www.wien.gv.at/gogv/19radabstellplaetze2016"/>
    </dcat:Distribution>
  </dcat:distribution>
</dcat:Dataset>
```

SPARQL endpoint - example



Retrieve all the resources from a dataset with a title that contains specific words (eg. 'Vienna')

SPARQL

You can search for the metadata stored in the EU Open Data Portal triple store by using the SPARQL endpoint query editor below.

Namespaces *

```
PREFIX dcat: <http://www.w3.org/ns/dcat#>
PREFIX odp: <http://data.europa.eu/euodp/ontologies/ec-odp#>
PREFIX dc: <http://purl.org/dc/terms/>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
```

SPARQL query *

```
SELECT ?DatasetTitle ?Publisher ?ResourceDescription WHERE { graph ?g {?DatasetURI a dcat:Dataset;
dc:publisher ?Publisher; dc:title ?DatasetTitle; dcat:distribution ?Resource. ?Resource dc:description ?
ResourceDescription. FILTER(regex(?DatasetTitle,"Vienna","i")) } } LIMIT 10
```

SPARQL endpoint - example



DatasetTitle	Publisher	ResourceDescription
"Glossary City of Vienna"	http://publications.europa.eu/resource/authority/corporate-body/CNECT	"Data archive containing files in the following formats: application/xml"
"University of Vienna Termbanks"	http://publications.europa.eu/resource/authority/corporate-body/CNECT	"Data archive containing files in the following formats: application/xml"
"Audioguide for the Military History Museum in Vienna"	http://publications.europa.eu/resource/authority/corporate-body/CNECT	"Data archive containing files in the following formats: MS-Word doc"

Schema.org

Started by Google, Microsoft, Yahoo, and Yandex to help with indexing web pages for search
Schema.org metadata can be embedded using microdata, RDFa or JSON-LD

Commonly used types

- Creative works: CreativeWork, Book, Movie, MusicRecording, Recipe, TVSeries
- Embedded non-text objects: AudioObject, ImageObject, VideoObject
- Event
- Organization
- Person
- Place, LocalBusiness, Restaurant ...
- Product, Offer, AggregateOffer
- Review, AggregateRating


schema.org

Schema.org

Google search results for "korona kielce". The search bar shows "korona kielce" and the results page displays several search results. A red box highlights a structured data snippet for "Korona Kielce" on the right side of the page.

Structured data

Korona Kielce
Football club



Korona Kielce, is a Polish football club, currently playing in the Ekstraklasa. In the years 2002–08 Club belonged to Polish holding company Kolporter Holding and achieved its greatest success – in 2005, winning promotion to the first division. Since then Korona has spent 5 seasons in Polish soccer top level. [Wikipedia](#)

Arena/Stadium: Kielce City Stadium
Location: Kielce, Poland
Managers: Maciej Bartoszek, Mirosław Smyła

Players

Jacek Kielb Midfielder	10
Petteri Forsell Midfielder	70
D'sean Theobalds Midfielder	26

[View 30+ more](#)

Google Dataset Search

schema.org:Dataset

- based on W3C DCAT

Full definition

- <https://schema.org/Dataset>

Google has an *application profile*

- <https://developers.google.com/search/docs/data-types/dataset>

Required properties	
description	<p>Text</p> <p>A short summary describing a dataset.</p>
name	<p>Text</p> <p>A descriptive name of a dataset. For example, "Snow depth in Northern Hemisphere".</p>

Recommended properties	
------------------------	--

Schema.org - example

View page source of any dataset at data.gv.at

- Navigate to <script> section

Search in Google Dataset Search for 'Radabstellplätze Wien 2016'

The image shows two overlapping screenshots. The left screenshot is from the data.gv.at website, displaying the dataset 'Radabstellplätze Wien 2016'. The right screenshot is from Google Dataset Search, showing the search results for the same dataset. A code block is overlaid on the screenshots, showing the JSON-LD schema for the dataset.

data.gv.at - Open Data Österreich

Katalog
Radabstellplätze Wien

Radabstellplätze in Wien

Radabstellplätze 2016 [CSV](#)

Daten und Ressourcen

Titel und Beschreibung	Number of bike storages
Englisch	
Veröffentlichende Stelle	Stadt Wien
Kontaktseite der veröffentlichenden Stelle	https://digitales.wien.gv.at
Veröffentlichende Stelle - E-Mailkontakt	open@post.wien.gv.at
Datenverantwortliche Stelle	Magistrat Wien - Magistratsabteilung 20 - Energieplanung
Kontaktseite der datenverantwortlichen Stelle	https://www.wien.gv.at/kontakte/ma20/index.html
Datenverantwortliche Stelle - E-Mailkontakt	post@ma20.wien.gv.at
Lizenz	Creative Commons Namensnennung 4.0 International
Lizenz Zitat	Datenquelle: Stadt Wien – https://data.wien.gv.at
Link zur Lizenz	https://creativecommons.org/licenses/by/4.0/deed.de
Link zu den Nutzungsbedingungen	https://data.wien.gv.at/nutzungsbedingungen
Attributbeschreibung	NUTS1 (z.B. AT1 für Ostösterreich) NUTS2 (z.B. AT13 für Bundesland Wien) NUTS3(z.B. AT130 für Stadt Wien) DISTRICT_CODE (z.B. 90001 für Wien) SUB_DISTRICT_CODE nicht verwendet) YEAR (Jahr, für das die Werte gelten) REF (Datenjahr) NUMBER (Anzahl der Radabstellplätze)
Geographische Abdeckung/Lage	Wien

```
<script type="application/ld+json">
{
  "@context": {
    "rdf": "http://www.w3.org/1999/02/22-rdf-syntax-ns#",
    "rdfs": "http://www.w3.org/2000/01/rdf-schema#",
    "schema": "http://schema.org/",
    "xsd": "http://www.w3.org/2001/XMLSchema#"
  },
  "@graph": [
    {
      "@id": "_:N7d7157714cf1492a944372bc2c224f3f",
      "@type": "schema:ContactPoint",
      "schema:contactType": "customer service",
      "schema:email": "open@post.wien.gv.at",
      "schema:name": "Magistrat Wien - Magistratsabteilung 20 - Energieplanung",
      "schema:url": "https://www.data.gv.at"
    },
    {
      "@id": "https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe",
      "@type": "schema:DataDownload",
      "schema:encodingFormat": "CSV",
      "schema:name": "Radabstellpl1u080e4tze 2016",
      "schema:url": "https://www.wien.gv.at/gogv/19radabstellplaetze2016"
    }
  ]
}
```

Google Dataset Search Radabstellplätze Wien 2016

1 result found

Radabstellplätze Wien
www.data.gv.at
Updated 01.03.2019
Dataset updated 01.03.2019
Dataset published 10.12.2018

Dataset provided by
Stadt Wien

License
<https://creativecommons.org/licenses/by/4.0/deed.de>

Available download formats from providers
CSV

Description
Radabstellplätze in Wien

Katalog Radabstellplätze Wien

Radabstellplätze in Wien

Daten und Ressourcen

Radabstellplätze 2016 **CSV**

[Mehr Information](#)

[Zur Ressource](#)

Titel und Beschreibung Englisch	Number of bike storages
Veröffentlichende Stelle	Stadt Wien
Kontaktseite der veröffentlichenden Stelle	https://digitales.wien.gv.at
Veröffentlichende Stelle - E- Mailkontakt	open@post.wien.gv.at
Datenverantwortliche Stelle	Magistrat Wien - Magistratsabteilung 20 - Energieplanung
Kontaktseite der datenverantwortlichen Stelle	https://www.wien.gv.at/kontakte/ma20/index.html
Datenverantwortliche Stelle - E-Mailkontakt	post@ma20.wien.gv.at
Lizenz	Creative Commons Namensnennung 4.0 International
Lizenz Zitat	Datenquelle: Stadt Wien – https://data.wien.gv.at
Link zur Lizenz	https://creativecommons.org/licenses/by/4.0/deed.de
Link zu den Nutzungsbedingungen	https://data.wien.gv.at/nutzungsbedingungen
Attributbeschreibung	NUTS1 (z.B. AT1 für Ostösterreich) NUTS2 (z.B. AT13 für Bundesland Wien) NUTS3(z.B. AT130 für Stadt Wien) DISTRICT_CODE (z.B. 90001 für Wien) SUB_DISTRICT_CODE(0 da nicht verwendet) YEAR (Jahr, für das die Werte gelten) REF_YEAR (Datenjahr) NUMBER (Anzahl der Radabstellplätze)
Geographische Abdeckung/Lage	Wien

```

<script type="application/ld+json">
{
"@context": {
  "rdf": "http://www.w3.org/1999/02/22-rdf-syntax-ns#",
  "rdfs": "http://www.w3.org/2000/01/rdf-schema#",
  "schema": "http://schema.org/",
  "xsd": "http://www.w3.org/2001/XMLSchema#"
},
"@graph": [
  {
    "@id": "_:N7d7157714cf1492a944372bc2c224f3f",
    "@type": "schema:ContactPoint",
    "schema:contactType": "customer service",
    "schema:email": "open@post.wien.gv.at ",
    "schema:name": "Magistrat Wien - Magistratsabteilung 20 - Energieplanung",
    "schema:url": "https://www.data.gv.at"
  },

```

```

    "@id": "https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe",
    "@type": "schema:Dataset",
    "schema:dateModified": "2019-03-01T10:20:42.981483",
    "schema:datePublished": "2018-12-10T14:46:53.400738",
    "schema:description": "Radabstellpl\u00e4tze in Wien\r\n",
    "schema:distribution": {
      "@id": "https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe/resour
    },
    "schema:includedInDataCatalog": {
      "@id": "https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe"
    }
  },

```

Google Dataset Search Beta

Search for Datasets



Try [boston education data](#) or [weather site:noaa.gov](#)

[Learn more](#) about including your datasets in Dataset Search.



1 result found



Radabstellplätze Wien

www.data.gv.at

Updated 01.03.2019

Radabstellplätze Wien



Dataset updated 01.03.2019

Dataset published 10.12.2018

Dataset provided by

Stadt Wien

License

<https://creativecommons.org/licenses/by/4.0/deed.de>

Available download formats from providers

CSV

Description

Radabstellplätze in Wien



Not seeing a result you expected?

[Learn](#) how you can add new datasets to our index.

Structured Data Testing Tool

The screenshot displays the Google Structured Data Testing Tool interface. The left pane shows the HTML source code of a dataset page, and the right pane shows the structured data extracted from the page.

Dataset

Property	Value
ID	https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe
@type	Dataset
@id	https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe
dateModified	2019-03-01T10:20:42
datePublished	2018-12-10T14:46:53
description	Radabstellplätze in Wien
keywords	fahrräder
keywords	verkehr
keywords	fahrrad
keywords	räder
license	https://creativecommons.org/licenses/by/4.0/deed.de
name	Radabstellplätze Wien
url	https://www.data.gv.at/katalog/dataset/radabsttelplstze-wien
distribution	
@type	DataDownload
@id	https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe/resource/e5dea2b6-7246-4617-9a28-4460dfc22e16
encodingFormat	CSV
name	Radabstellplätze 2016
uri	https://www.wien.gv.at/govy/l3radabsttelplstze2016
includedInDataCatalog	
@type	DataCatalog
@id	https://www.data.gv.at/katalog/dataset/_N445cee4dae9c45488bce0f098c998f16

<https://search.google.com/structured-data/testing-tool>

Structured Data Testing Tool

Dataset		0 ERRORS 0 WARNINGS ^
ID: https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe		
@type	Dataset	
@id	https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe	
dateModified	2019-03-01T10:20:42	
datePublished	2018-12-10T14:46:53	
description	Radabstellplätze in Wien	
keywords	fahrräder	
keywords	verkehr	
keywords	fahrrad	
keywords	räder	
license	https://creativecommons.org/licenses/by/4.0/deed.de	
name	Radabstellplätze Wien	
url	https://www.data.gv.at/katalog/dataset/radabstellplatze-wien	
distribution		
@type	DataDownload	
@id	https://www.data.gv.at/dataset/af8e02b6-1e03-4464-a69b-8533d8703ffe/resource/e5dea2b6-7246-4617-9a28-4460dfc22e16	
encodingFormat	CSV	
name	Radabstellplätze 2016	
url	https://www.wien.gv.at/gogv/l9radabstellplaetze2016	

Synergy between DOIs and schema.org

DataCite provides metadata on DOIs

- RDFs compliant with schema.org

Example



- <https://data.datacite.org/application/vnd.schemaorg.ld+json/<DOI>>
- <https://data.datacite.org/application/vnd.schemaorg.ld+json/10.1371/journal.pcbi.1006750>

JSON-LD returned by DataCite can be directly embedded on web pages

- Performance issues!



<https://data.datacite.org>



Repository registries

Directory of Open Access Repositories – DOAR

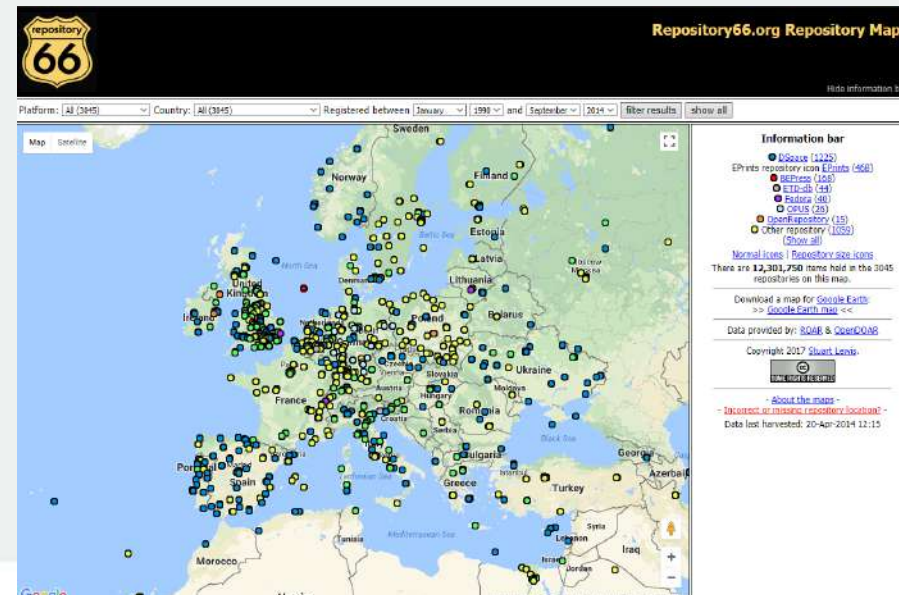
- Based on registrations
- <http://www.opendoar.org/>

Registry of Open Access Repositories – ROAR


- Automatically harvested list based on OAI-PMH
- <http://roar.eprints.org/>

Projection of DOAR and ROAR on maps

- <http://maps.repository66.org>
- re3data.org
- FAIRsharing.org



Repository registries – re3data

Repository details 

Phaidra Universität Wien

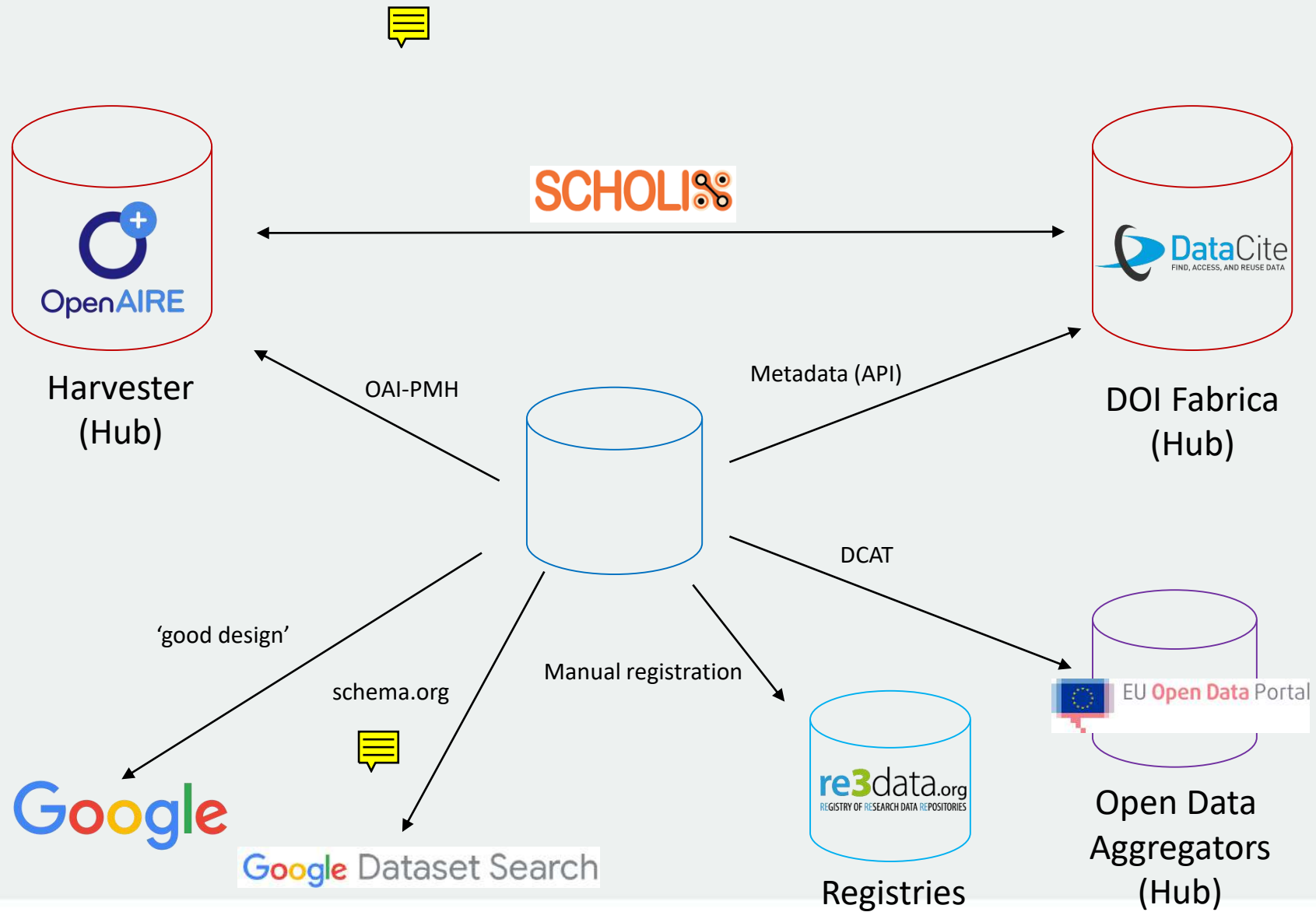
General Institutions Terms Standards

Name of repository	Phaidra Universität Wien
Additional name(s)	Permanent Hosting, Archiving and Indexing of Digital Resources and Assets
Repository URL	https://phaidra.univie.ac.at/
Subject(s)	Humanities and Social Sciences Life Sciences Natural Sciences Engineering Sciences
Description	Phaidra Universität Wien, is the innovative whole-university digital asset management system with long-term archiving functions, offers the possibility to archive valuable data university-wide with permanent security and systematic input, offering multilingual access using metadata (data about data), thus providing worldwide availability around the clock. As a constant data pool for administration, research and teaching, resources can be used flexibly, where continual citability allows the exact location and retrieval of prepared digital objects.
Contact	support.phaidra@univie.ac.at phaidra@univie.ac.at
Content type(s)	Images Audiovisual data Scientific and statistical data formats Networkbased data Plain text other
Keyword(s)	hosting long-term-archiving multidisciplinary digital objects research
Repository type(s)	institutional other
Mission statement for designated community	https://datamanagement.univie.ac.at/en/about-phaidra/policy-of-phaidra/
Research data repository language(s)	eng deu ita
Data and/or service provider	dataProvider

<https://www.re3data.org/repository/r3d100010472>

SUMMARY

External visibility – your main options



You should know

How to make repository contents visible?

- what options there are and how to choose the best one for you setting?
- how to describe data using discussed standards
 - DCAT
 - DataCite

How do repositories support FAIRness?

- Which (sub-) principles specifically and how?

Next lecture

Developing Research Data Management Services

- Data lifecycle model
- Policies
- Costs and Business Models
- Data Stewards
- Repository Certification

Developing Research Data Management Services and Repository Certification

Tomasz Miksa

Agenda

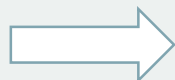
Research Lifecycle models

Developing Research Data Management (RDM) Services

- Policies
- Costs
- Data Management
- Infrastructure

Repository Certification

Introduction




Perspective change



- 'This lecture
 - You need to design a solution at an institution supporting data management and preservation
 - Taking into account technical, organisational, cultural, political problems...
 - **What do you do?!**

The size and governance structure of institutions has an impact

Large, hierarchical institutions

- Move slowly
- Require a lot of advocacy 
- Have more resource
- Economies of scale



Smaller institutions

- More agile
- Simpler communication
- More focussed vision
- Less resource



Establish a long term governance group

Good mix of representatives from operational units

Senior management leadership

Avoid relying on one or two key individuals

Keep the group active

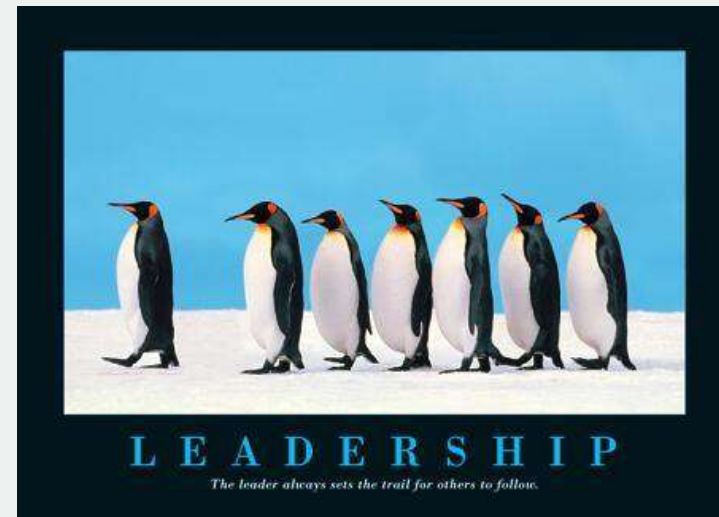



Image credit: <http://www.executive-coaching-services.co.uk/executive-coaching/leaders.jpg>

Be aware of existing infrastructure

 In many HEIs services are beginning to be embedded but aren't joined up effectively

Research Organisation Description		
RDM profile component	Record HEI Link - insert your URLs in the space provided	Guidance
Means of raising staff awareness of funders' research data requirements		Provide a link to an information page on funders' policies. This could be internal or external (E.g., DCC's policy overview table)
Research data policy		Provide a link to research data policy or aspirational statement
Strategy or implementation plan for research data services		Provide a link to research data strategy page or roadmap
RDM advice and support services		Provide a link to page describing data management planning guidance and/or support services at this organisation
Active data storage		Provide link(s) to active research data storage information page(s). There may be multiple options at Research Group/School/College/Central levels.
Data register or catalogue		Provide a link to your internal research data registration homepage. This may be provided via the data repository and/or CRIS.
Persistent identification for datasets		Provide a link to any page(s) detailing schemes used to identify digital data items (e.g., DataCite).
Data access procedures		Provide a link to any information provided about research data access.
Secure data access		Provide a link to any information provided about secure data access and governance.
Institutional publications repository (if it includes research data or metadata)		Provide a link to your institutional repository homepage
Data repository for longer term access and preservation		Provide a link to your research data repository homepage. This may be an extension of your publications repository, a separate data repository or a pointer to an external data repository service (E.g., Zenodo).

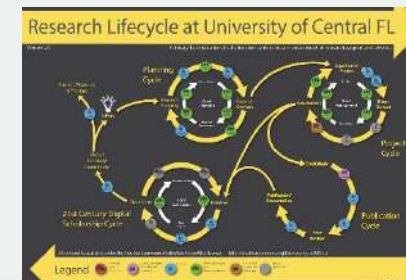
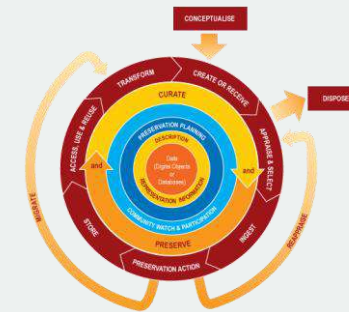
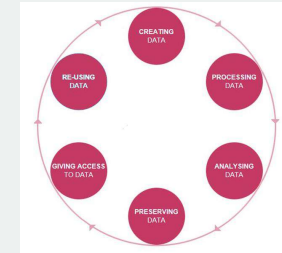


<http://www.dcc.ac.uk/projects/opd-for-rdm>

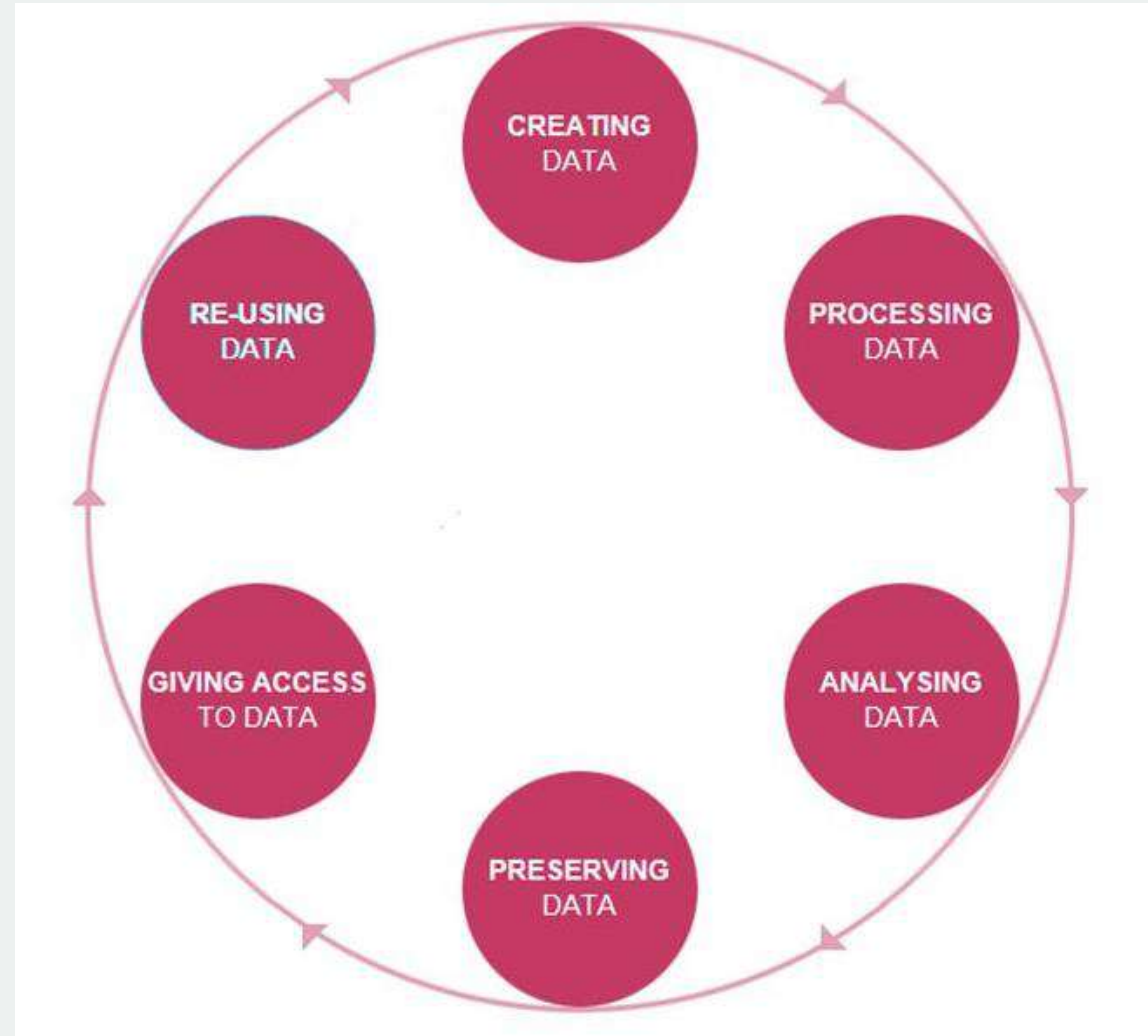
Understand processes

Research Data Lifecycle Models

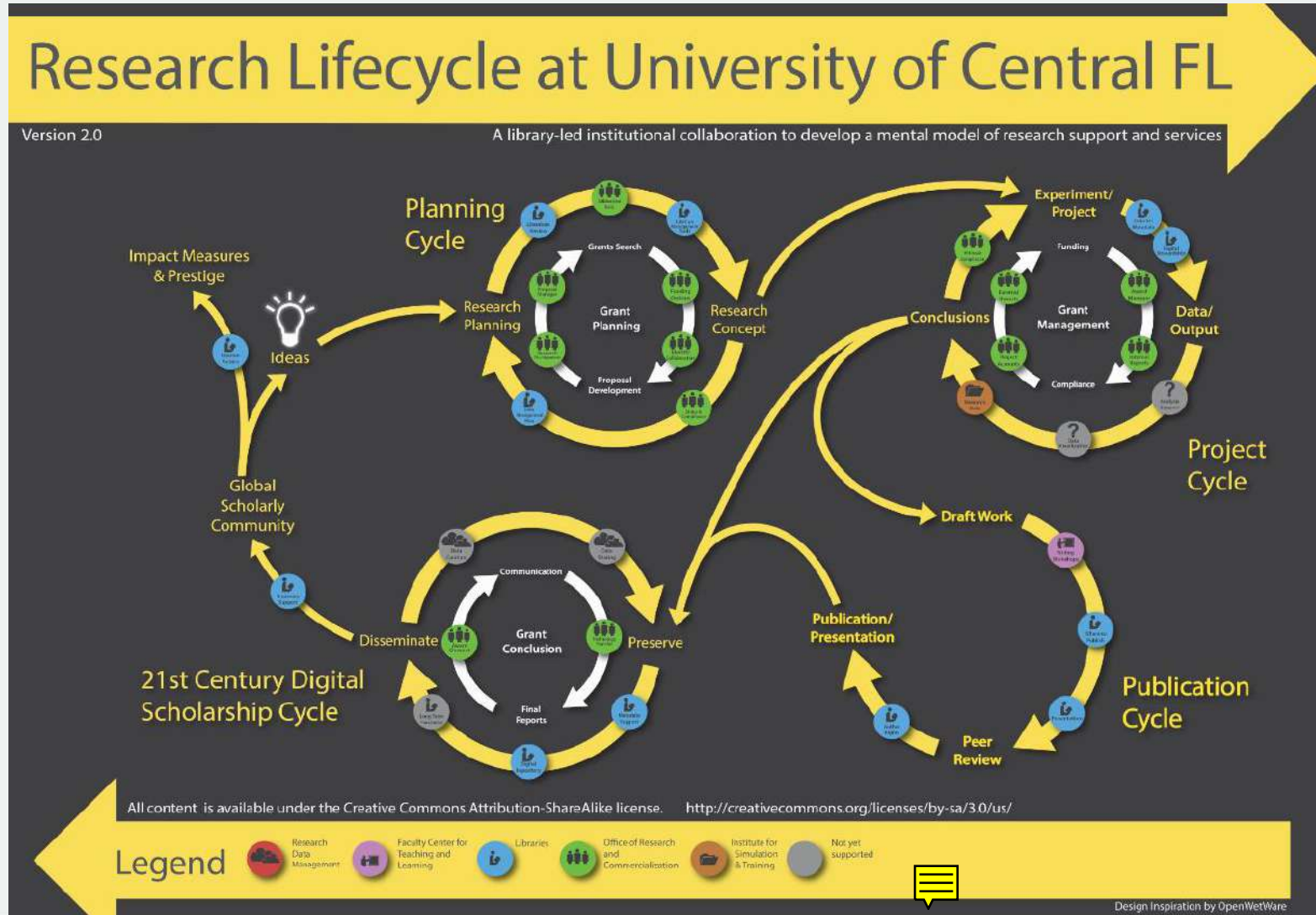
- Describe roles of stakeholders
- Help in
 - tailoring services
 - identifying responsibilities
 - defining infrastructure
- NOT to be used by researchers 
- Examples
 - UK Data Archive
 - Digital Curation Centre
 - University of Central Florida



UK Data Archive Lifecycle model



University of Central Florida Lifecycle Model

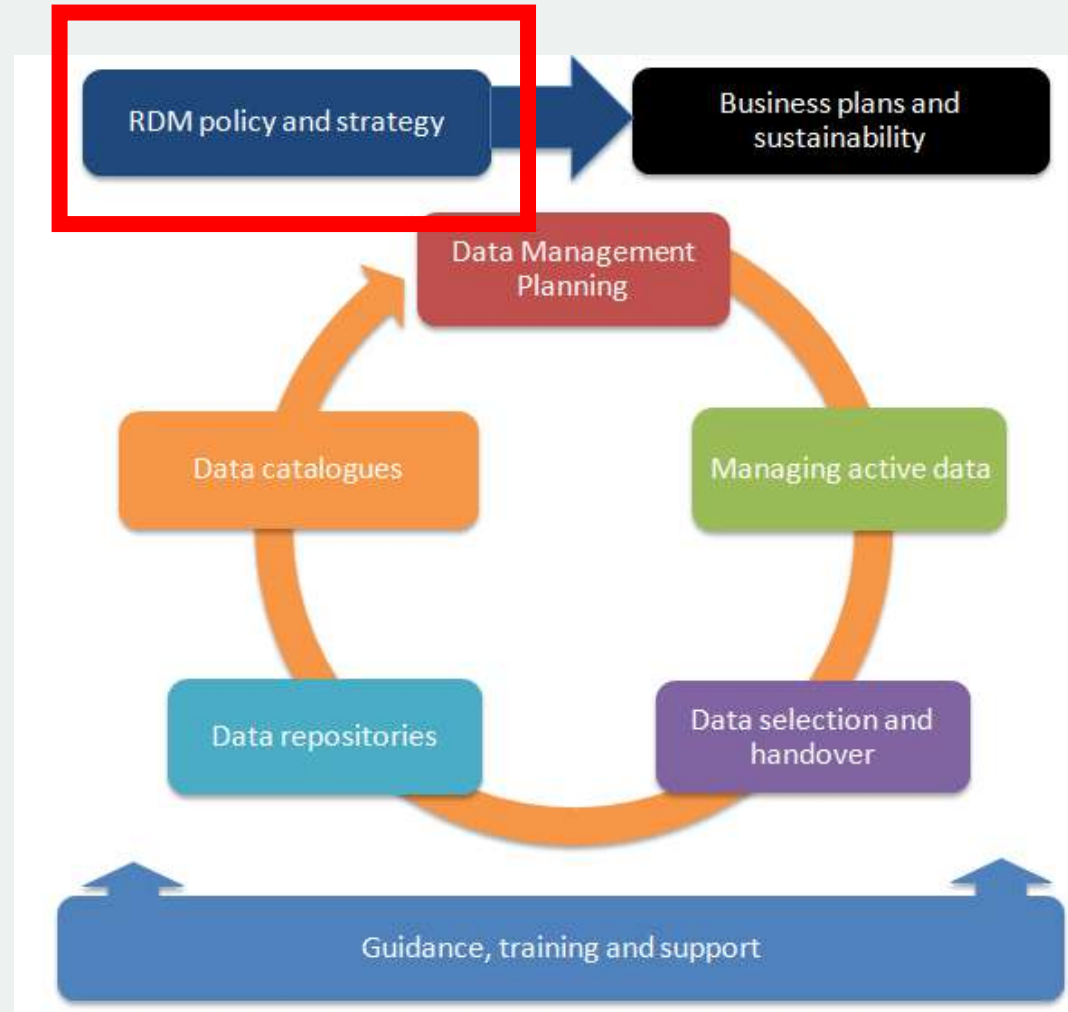


Developing Research Data Management (RDM) Services

Components of an RDM service infrastructure

- RDM policy and strategy
- Business plans and sustainability
- Guidance, training and support
- Data management planning
- Managing active data
- Data selection and handover
- Data repositories
- Data catalogues

Developing Research Data Management (RDM) Services



Understanding policies



1

Starting with some *Taboos*

2

Going over to related *Principles*

3

Going over to the creation of a *Policy*

4

Going over to *Rules, Legislations* and *Regulations* (canons, norms, guidelines)

Taboo

A **taboo** is something, which is **forbidden** or **disapproved of**, or placed under a social prohibition.

“Thou shalt not delete scientific data“


“Thou shalt not destroy infrastructures”

Usually a **negative** assertion.

In society and academic environment taboos are accepted only if they are just **a few**.



Principle

A principle is a fundamental truth or proposition that serves as the foundation for a system of belief 

*Research **data** are to be preserved*

*Research **data** are to be kept **FAIR** - Findable, Accessible, Interoperable, Reusable.*

*Research data **infrastructures** are to be kept accessible*

Format: positive assertion

Policies

A policy is...

- a **course of action** or **principles** adopted by an organization

“The Institution [name XY] will preserve its research data infrastructure always accessible and free to its members according to the FAIR principles”

General assumptions concerning policies:

- A single Policy
 - the policy is a single entity, it should not be in competition with other policies
- Creators of Policy do not want to modify it
- Policy is usually accepted after a while
- Policy offers the frame for the generation of Rules




Rules, Regulations

Rules are prescribing conducts or actions

They are generated by the founder of “orders”

Characteristics of rules are:

- There may exist “**lots of rules**”
- Rules are not always **clear**
(they often need interpretation according to the situation)
- Rules are **usually accepted**, but **often imposed** procedures
- It is allowed to **modify Rules**
- The Law is an expression of rules 

Rules, Regulations

Example:

“Our University will maintain accessible our infrastructure each day from 9:00 a.m. to 12:00 a.m and offer support only on Friday from 7:00 a.m. to 8:00 a.m. The research data, that are publicly funded are to be kept free and accessible to all members of our University each Sunday, from 9:00 to 12:00 a.m.”

Taboos	Principles	Policies	Rules
<p>Negative assertion</p> <p>few</p> <p><i>“You shall not delete scientific data”</i></p> <p><i>“You shall not destroy infrastructures”</i></p>	<p>Positive assertion</p> <p>more than „few“</p> <p><i>“Research data are to be kept FAIR - Findable, Accessible, Interoperable, Reusable.”</i></p> <p><i>“Research data infrastructures are to be kept accessible”</i></p>	<p>A course or principle of action. Policy offers the frame for the generation of Rules, should not be in competition with other policies</p> <p><i>“The Institution [name XY] will preserve its research data infrastructure always accessible and free to its members according to the FAIR principles”</i></p>	<p>Rules prescribe conducts or actions; define who what when and where should be done according to the Policy</p> <p><i>“Our University will maintain accessible our infrastructure each day from 9:00 a.m. to 12:00 a.m and offer support only on Friday from 7:00 a.m. to 8:00 a.m. “</i></p>

Why these differentiations?

It is important to identify the different semantic levels

Understanding of the semantic hierarchy is useful in order to produce appropriate guidelines

Policy – LEARN project

Leader Activating Research Networks Research Data Management Policy

- can be tailored by any University or Research Institution
- based on existing European policies



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 654139.

Outreach example: Austria



Adaptation to needs of

- five Austrian art universities (started)
- three Medical Universities
- TU Wien

e-infrastructures
austria

Policy sections

1. Preamble
2. Jurisdiction
3. Intellectual Property Rights
4. Handling Research Data
5. Responsibilities, Rights, Duties
 - 5.1 Researchers are responsible for:...
 - 5.2 The [name of research institution] is responsible for:
6. Validity



<http://learn-rdm.eu/wp-content/uploads/RDMToolkit.pdf?pdf=RDMToolkit>

Policy section examples

5.1 Researchers are responsible for:

- Collecting, documenting, archiving, providing access to and storing or ensuring the proper destruction of research data and records. (...) Such information should be included in a Data Management Plan (DMP) that explicitly defines the approach to matters of data collection, administration, integrity, confidentiality, storage, use and publication.

5.2 The Institution is responsible for:

- Designing and deploying mechanisms and services for the storage, safekeeping, registration and deposition of research data to support current and future access to research data during and after the completion of research activities;

Tricky points:

- How do you define research data? (publications, source code, raw data..)
- Who are researchers? (students, employed formally, etc.)

Homework – read TUW policy

Research > RTI support > Research data > Research data management > Policy

Policy for Research Data Management (RDM) at the TU Wien

The TU Wien sees itself as playing an important role in the expansion of knowledge and technology transfer of research results and thus encouraging innovation and ultimately benefitting economy and society. A key to supporting academic research activities lies in the institution's ability to systematically manage, preserve and make available scientific output from different disciplines for reuse. In its Policy for Research Data Management, the TU Wien affirms the value of research data for research and teaching and the potential of their reuse by society.

↓ Policy_for_Research_Data_Management.pdf PDF 548 KB

1. PREAMBLE

2. SCOPE

3. RIGHTS OF USE

4. RESEARCH DATA AND RESEARCH DATA MANAGEMENT

5. HANDLING RESEARCH DATA

6. RESPONSIBILITIES, RIGHTS AND DUTIES

7. VALIDITY

Definitions

Recommendations

RDM policy and strategy

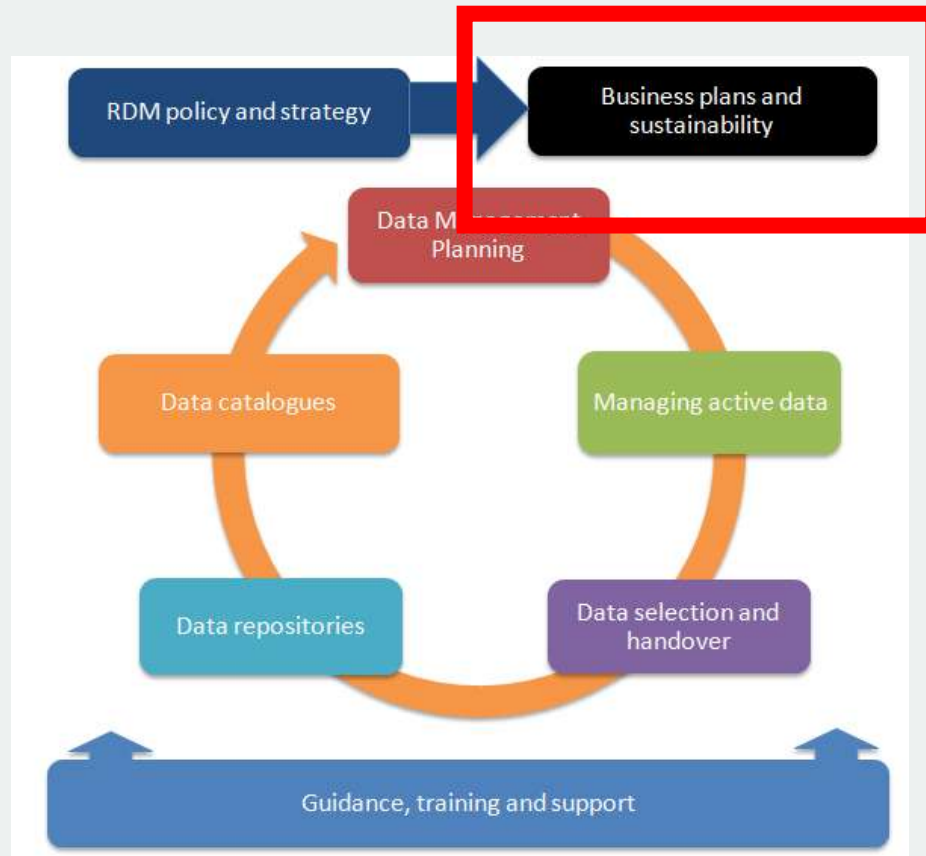
Develop a strategy

- Understand your current position and where you want to be to define your strategy
- Map out a programme of activity to deliver infrastructure and services

Develop a policy

- Draft a policy based on external drivers and local context to establish your core RDM principles
- Ratify the policy then undertake advocacy work and pilot studies to aid implementation
- Consult broadly to gain consensus and secure support





BUSINESS PLANS AND SUSTAINABILITY

Who pays for what?

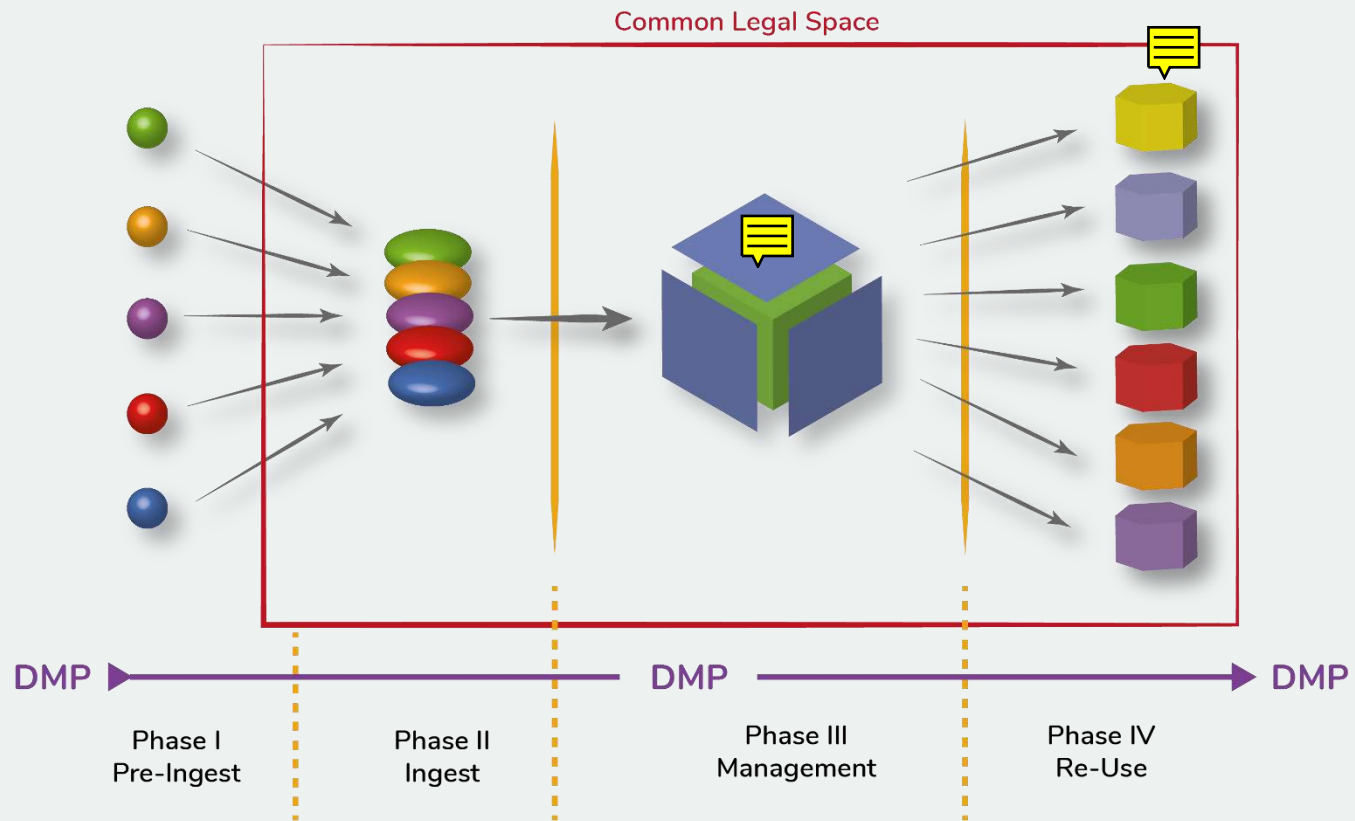


Data from various sources

Data conversion and enrichment

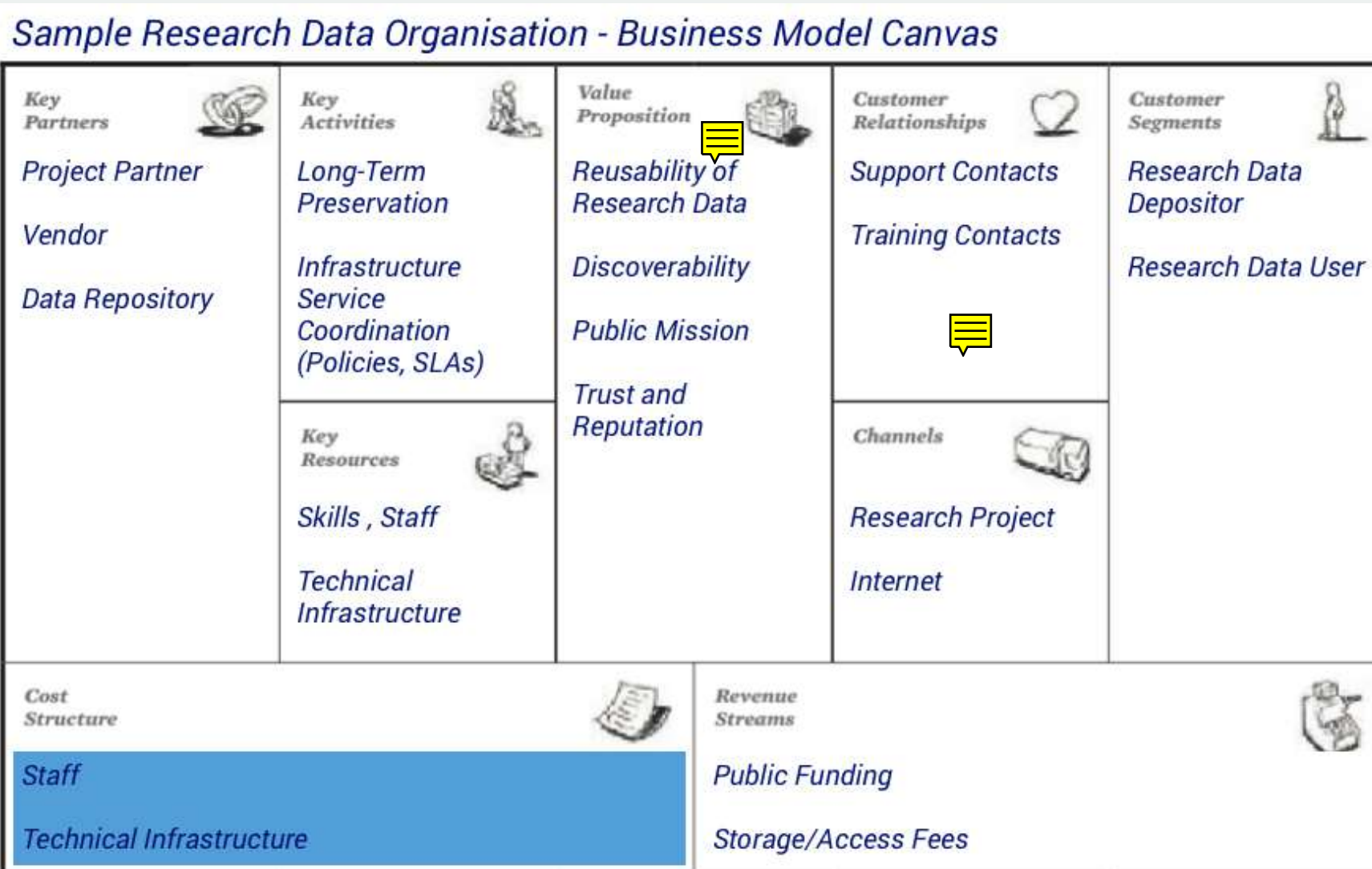
Data Preservation

Data Re-use



Digital workflow model for data management
Raman Ganguly, University of Vienna 2017

Identify the cost structure



Business plans and sustainability

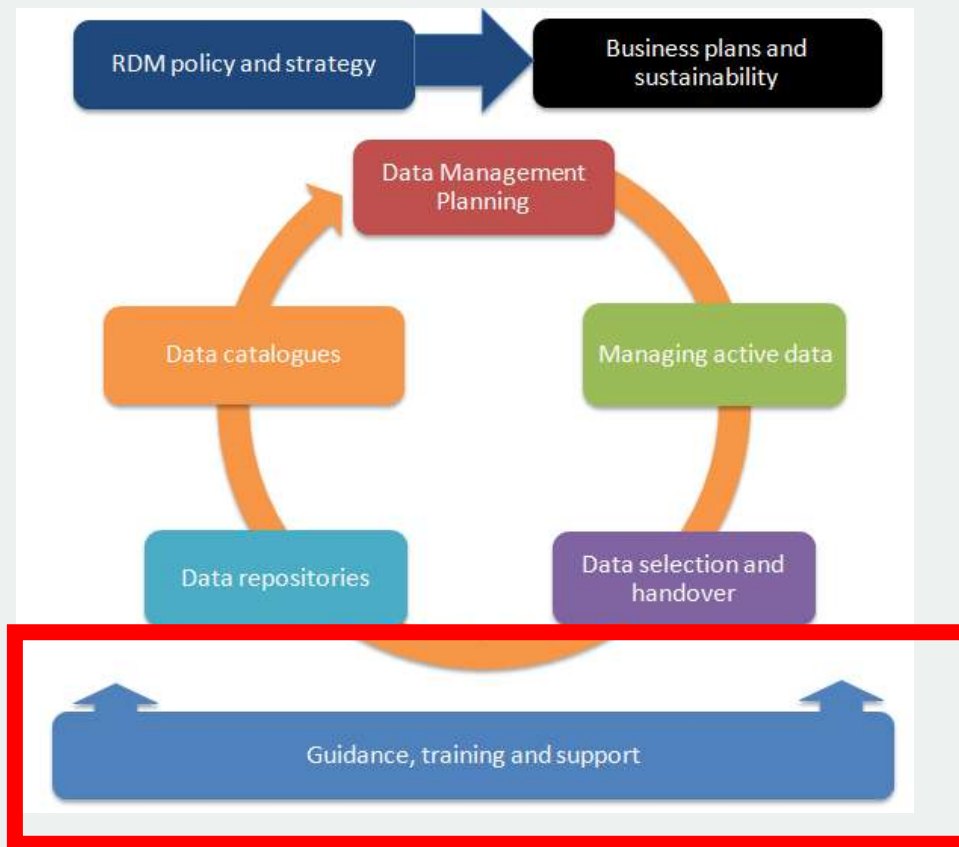
Based on your strategy, develop a phased **business plan** covering projections for 3, 5 and 10 years

Identify predicted **costs** and planned expenditures, together with an indication of the financial year in which expenditure will occur

Consider whether any costs can be recouped, for example by **charging for services**

Undertake a **cost/benefits analysis** to help make the case for investment

Address **sustainability** issues and the associated **long-term costs**



GUIDANCE, TRAINING AND SUPPORT

Guidance, training and support

Single Point of Information

- Website
 - Collate details of existing institutional support to provide a single, joined-up RDM guidance website
- Helpdesk
 - Coordinate the provision of support, either through a central helpdesk or well-signposted contacts

Trainings and Consultations

- Provide RDM training to support researchers and reskill support staff
- Offer more in-depth consultancy services



Singe point of information - example

RESEARCH DATA

← RTI SUPPORT

Overview

Research data management >

DMP >

Funders' guidelines >

Storing and sharing >

Preserving and publishing >

Persistent Identifiers >

Team

Contact

Search

TU WIEN STUDIES RESEARCH PARTNERSHIPS SERVICES

Welcome to the Center for Research Data Management

» Are you looking for suitable storage solutions for your data collection? Would you like to preserve your data, code and associated software and share it with others? Do you need a repository? Must you calculate data management costs for your project application, or are you required to create a data management plan? Has this raised legal or ethical questions? «

Team Center for Research Data Management

Then you have come to the right place. We, a team of TU Wien employees from various fields, offer you, a researcher and research group leader at the TU Wien, support in handling research data throughout the entire data life cycle.

Together with other departments at the TU Wien, we are contributing to a coordinated expansion of research data management services. We strive to adapt these services to your needs and constantly develop. Your feedback and suggestions are therefore welcomed.

What can we do for you?

research.data@tuwien.ac.at
T: +43-1-58801-44002
Give us a call, write to us or use our [online form](#).

To keep you informed about news, services and events, sign up for our [mailing list](#).

RESEARCH DATA MANAGEMENT

Research data management refers to the entire life cycle of research data.

DATA MANAGEMENT PLANS

A data management plan helps to make data accessible and reusable in the long term.

FUNDERS' GUIDELINES

Funding bodies are more and more interested in the research on which publications are based.

STORE AND SHARE DATA

Using the services of the TU Wien to store and share data is a prudent decision.

PRESERVE AND PUBLISH DATA

We offer practical alternatives to storing hard drives in your desk drawer.

PERSISTENT IDENTIFIERS

ORCID, DOI and co. distinguish you and your data.

<https://www.tuwien.at/en/research/rti-support/research-data/overview/>

Guidance, training and support

Who can do this?

What is the profile of the best candidate?

‘We need 500.000 respected data stewards to operate the European Open Science Cloud’

<http://e-irg.eu/news-blog/-/blogs/we-need-500-000-respected-data-stewards-to-operate-the-european-open-science-cloud>

Data stewards



Data Stewardship

TU Delft in NL as an example

- Subject-specific Data Steward at every Faculty
- Half-time positions
- PhD holders in a given area proffered
- Strategic coordination from the Library



What do data stewards do?

Analyse data management needs

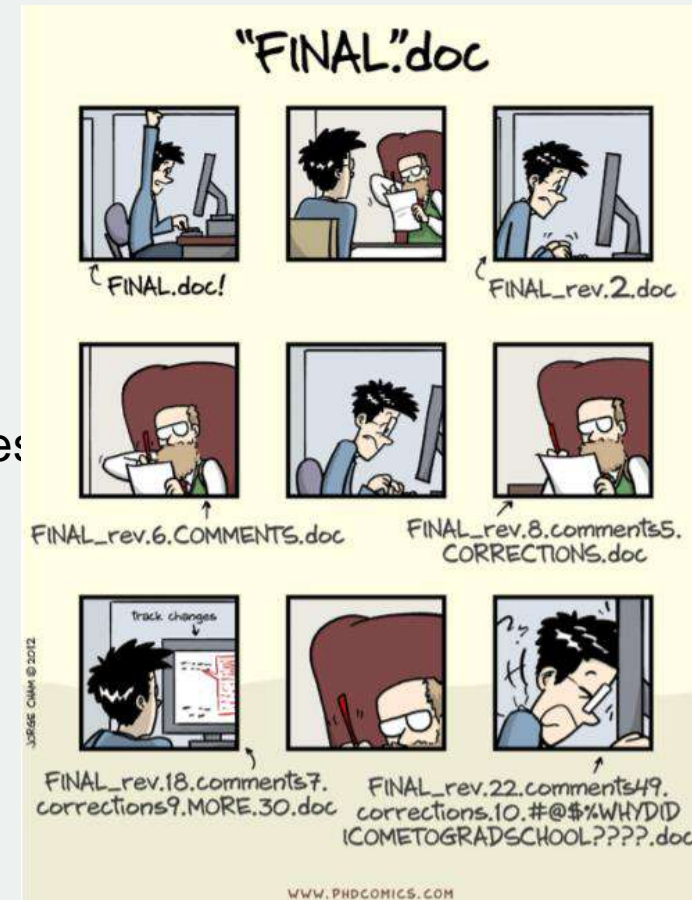
Provide advice to researchers

Train and inspire

Help comply with funders' and journals' policies

Develop faculty research data policies: roles and responsibilities

Trusted point of contact for data management questions



FEBRUARY 23, 2018

We are hiring (again!) – Data Steward position at TU Delft

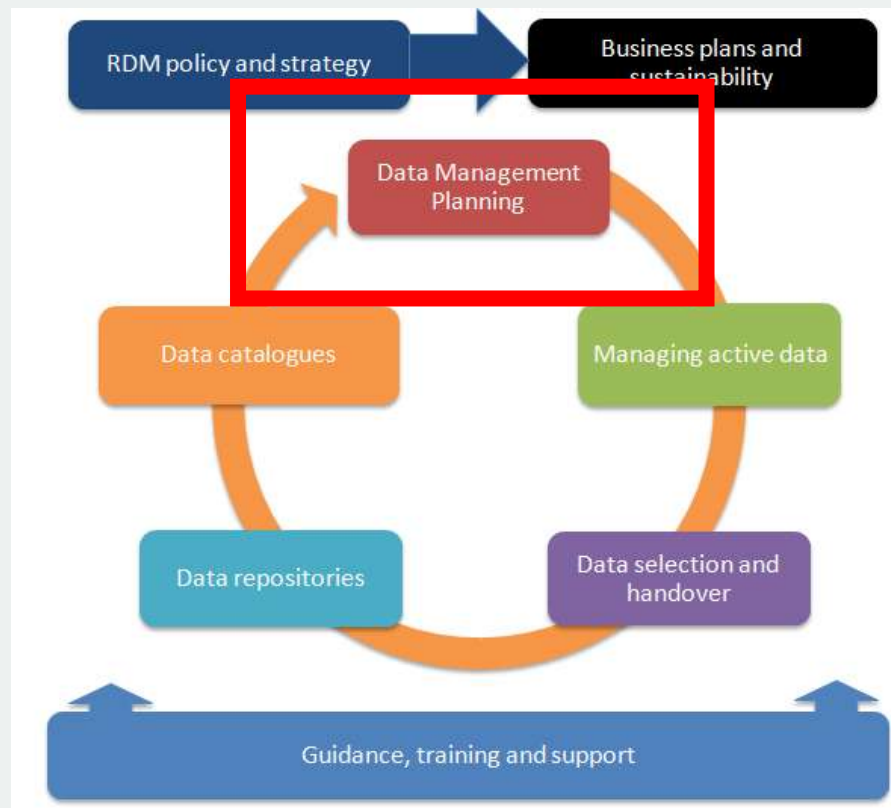


We have an exciting job opening for a Data Steward at TU Delft at the Faculty of Architecture & Built Environment and the Faculty of Industrial Design (joint appointment): <https://www.academictransfer.com/employer/TUD/vacancy/45483/lang/en/>

- Closing date: 15 March 2018
- Salary: up to € 4084/month
- We are looking for individuals enthusiastic about data management and who have a PhD degree in the relevant subject area (or equivalent experience).

This is a great chance to join the dynamically growing team of [Data Stewards at TU Delft](#) and to contribute to a cultural change in research data management in a disciplinary manner. The job is really about inspiring the research community and improving day to day practices, and not about policy compliance.

All informal inquiries can be directed to me: M.Teperek@tudelft.nl



DATA MANAGEMENT PLANNING

Data management planning

Data Management Plan (DMP) describes:

- how the data will be created
- how it will be documented
- who will be able to access it
- where it will be stored
- who will back it up
- whether (and how) it will be shared & preserved



Do you really need your own DMP template?

1. Does your institution encourage/require that researchers write Data Management Plans, even if their funder does not require one? If so, what information are researchers asked to provide?
2. Do DMPs submitted to research funders supply the information your institution asks for (if any), or are there additional questions that you want to ask?
3. What guidance, examples and suggested answers can you provide to help researchers write DMPs?

[<http://www.dcc.ac.uk/sites/default/files/documents/tools/dmpOnline/DMPonline-customisation-guidelines.pdf>]

If you need one...

Identify who and why will need this information

Browse existing checklists for help

Provide a template

- questions
- guidance
- possible answers

Identify overlaps of your template (e.g. with Horizon 2020)

Test it and incorporate feedback

It is not only about DMPs...

Design the whole data management process

- do you need different DMP phases?
- who receives the DMPs?
- who reviews them?
- who provides assistance?
- etc.

Provide a reviewers' guide

- clear judgment criteria

Offer default data location

- if you require people to deposit data and describe in a DMP, provide a place for doing that

It is not only about DMPs...

Establish a policy

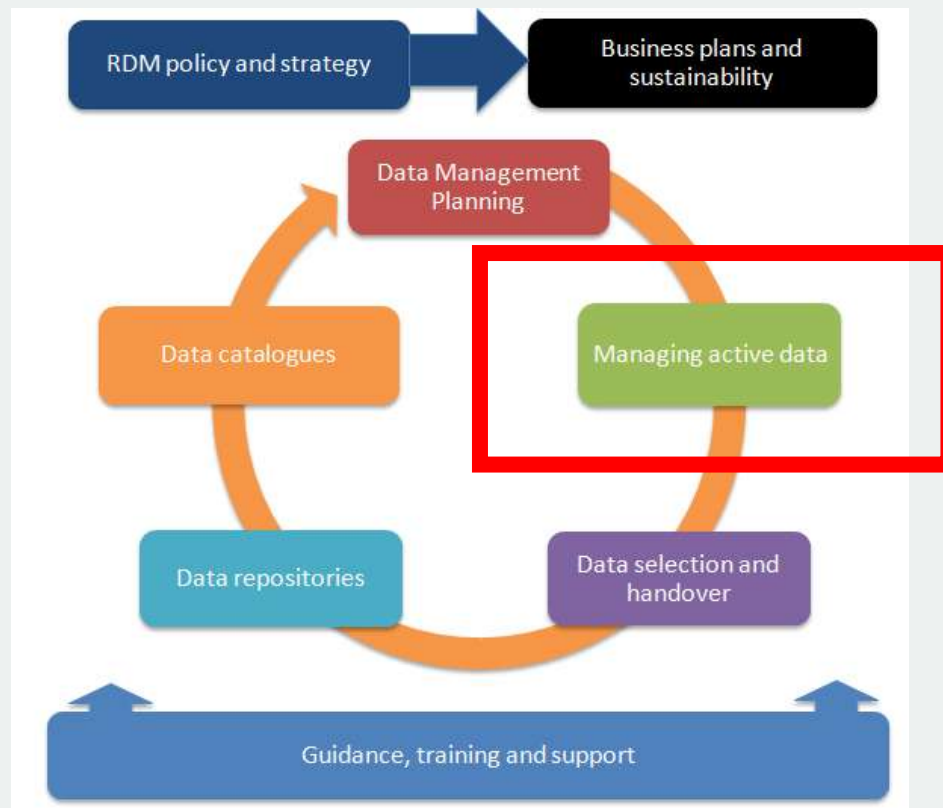
- make a use of DMPs obligatory
- consider incentives and penalties

Provide support

- tools
- examples
- trained staff for assistance in DMP creation
 - related problem: Where do you find people with expertise?

Awareness and understanding

- offer training material
- organise workshops
- related problem: How to teach correct data management practices at schools...



MANAGING ACTIVE DATA

Managing active data

Host on your own



Open Science Framework



GitHub



Outsource?

- Sensitive data?
- Jurisdiction?



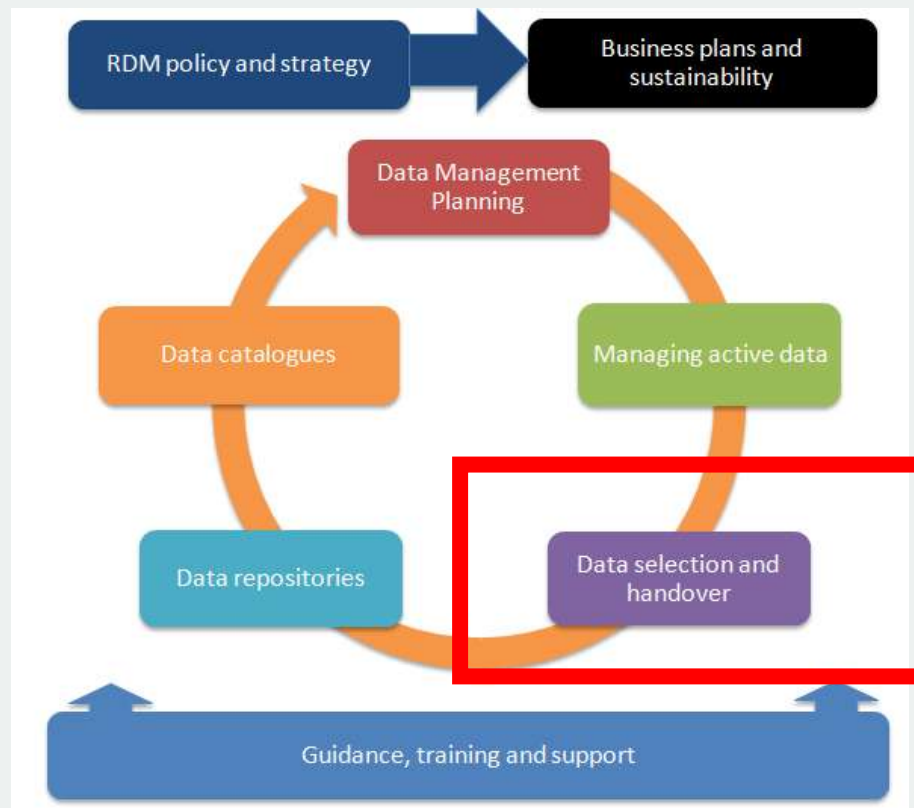
Managing active data

Review data holdings and RDM practices to see if the current infrastructure and systems are sufficient

Where appropriate, make a case for investment to provide additional research data storage

Develop procedures for the allocation and management of research data storage

Provide flexible RDM systems to support the creation, management and sharing of data that meet a diverse range of research contexts and needs



DATA SELECTION AND HANDOVER

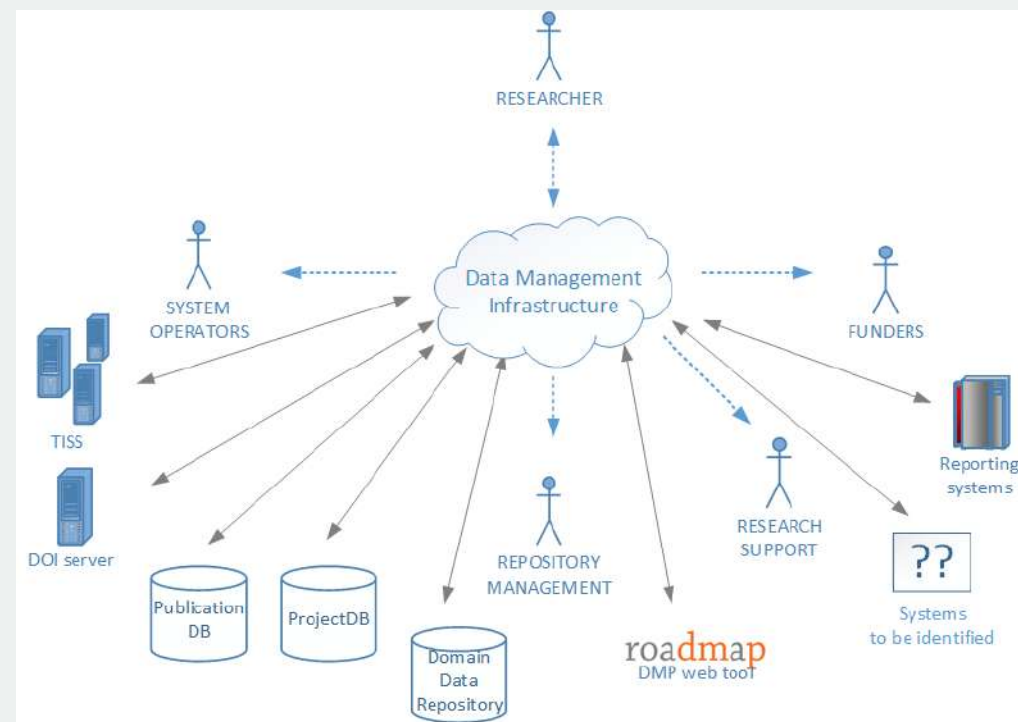
Data selection and handover

What needs to be kept?

- Do you encourage/require specific data to be deposited?
 - Into your own repository?

Data Management Infrastructure

- Integrate existing systems
- Seamless data flow



Data selection and handover

Identify which data fall under your remit and establish criteria to guide decisions on what to keep

If establishing a repository, develop deposit agreements and high-level guidance

Develop or use existing deposit tools to ease the process of handover

Advocate the benefits of data deposit
to encourage uptake

Support research groups to develop
guidance and offer input to decisions



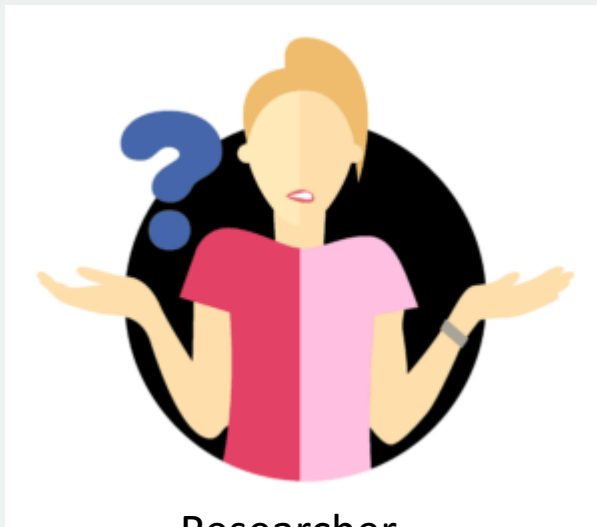


DATA REPOSITORIES

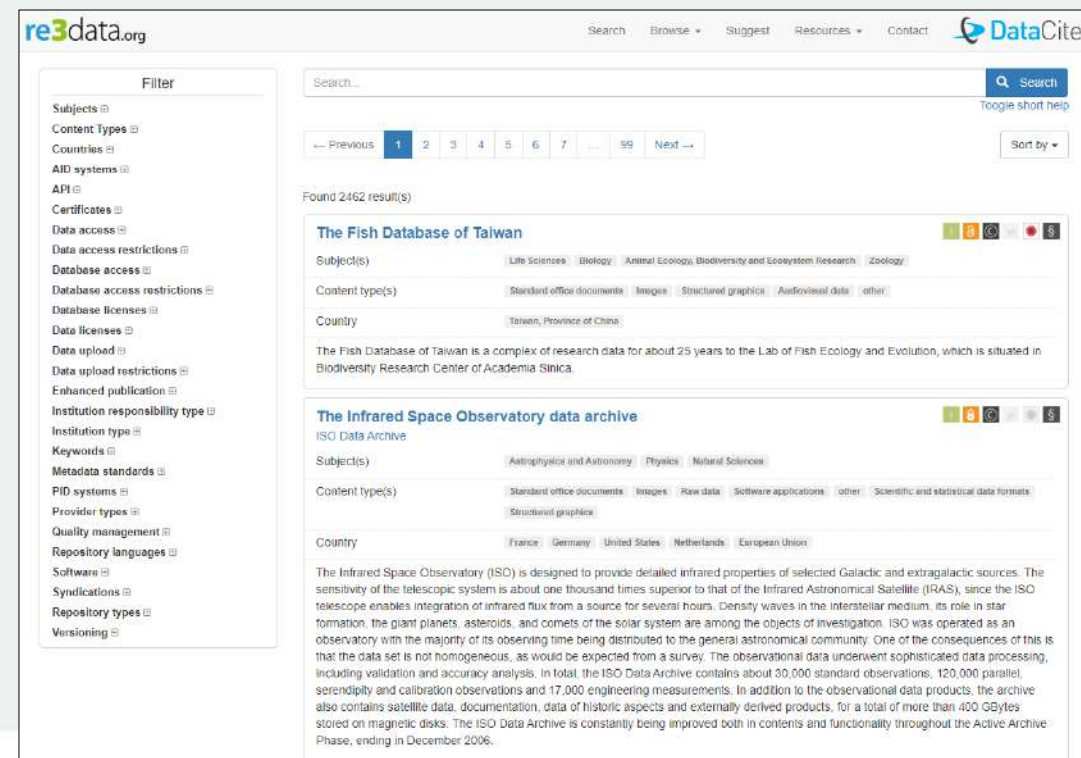
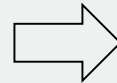
Data repositories

Which repository to use?

What are all these filters for?



Researcher



The screenshot shows the re3data.org website interface. At the top, there is a search bar and navigation links for Search, Browse, Suggest, Resources, and Contact. The main content area displays search results for 2462 items. Two results are visible:

- The Fish Database of Taiwan**: A database access repository. Subject(s): Life Sciences, Biology, Animal Ecology, Biodiversity and Ecosystem research, Zoology. Content type(s): Standard office documents, Images, Structured graphics, Audiovisual data, other. Country: Taiwan, Province of China. Description: The Fish Database of Taiwan is a complex of research data for about 25 years to the Lab of Fish Ecology and Evolution, which is situated in Biodiversity Research Center of Academia Sinica.
- The Infrared Space Observatory data archive**: An ISO Data Archive. Subject(s): Astrophysics and Astronomy, Physics, Natural Sciences. Content type(s): Standard office documents, Images, Raw data, Software applications, other, Scientific and statistical data formats. Country: France, Germany, United States, Netherlands, European Union. Description: The Infrared Space Observatory (ISO) is designed to provide detailed infrared properties of selected Galactic and extragalactic sources. The sensitivity of the telescopic system is about one thousand times superior to that of the Infrared Astronomical Satellite (IRAS), since the ISO telescope enables integration of infrared flux from a source for several hours. Density waves in the interstellar medium, its role in star formation, the giant planets, asteroids, and comets of the solar system are among the objects of investigation. ISO was operated as an observatory with the majority of its observing time being distributed to the general astronomical community. One of the consequences of this is that the data set is not homogeneous, as would be expected from a survey. The observational data underwent sophisticated data processing, including validation and accuracy analysis. In total, the ISO Data Archive contains about 30,000 standard observations, 120,000 parallel serendipity and calibration observations and 17,000 engineering measurements. In addition to the observational data products, the archive also contains satellite data, documentation, data of historic aspects and externally derived products, for a total of more than 400 GBytes stored on magnetic disks. The ISO Data Archive is constantly being improved both in contents and functionality throughout the Active-Archive Phase, ending in December 2006.

Data repositories

Research support

- Often has too little information on domain, data and needs

Not every researcher can get support



Data repositories: goal

Automate repository recommendation

- Set relevant filters automatically

Reduce effort

- Narrow down selection to relevant repositories

Lower the expertise needed

- Why should researchers deal with licenses if the funder/institutional policy regulates this anyway?

Automation example – RDM policy

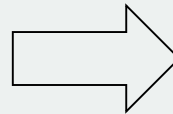
‘... research data should be assigned an open use license’

Policy for Research Data Management (RDM) at the TU Wien

The TU Wien sees itself as playing an important role in the expansion of knowledge and technology transfer of research results. It thus encourages innovation and ultimately benefits economy and society. A key to supporting academic research activities lies in the institution's ability to systematically manage, preserve and make available scientific output from different disciplines for reuse. In its Policy for Research Data Management, the TU Wien affirms the value of research data for research and teaching and the potential of their reuse by society.

Policy_for_Research_Data_Management.pdf PDF 546 KB

- 1. PREAMBLE
- 2. SCOPE
- 3. RIGHTS OF USE
- 4. RESEARCH DATA AND RESEARCH DATA MANAGEMENT
- 5. HANDLING RESEARCH DATA
- 6. RESPONSIBILITIES, RIGHTS AND DUTIES
- 7. VALIDITY
- Definitions
- Recommendations



re3data.org

Database licenses ⊕

Data licenses ⊖

- Apache License 2.0 (13)
- BSD (12)
- CC (97)**
- CC0 (97)**
- Copyrights (18)
- ODC (19)
- OGL (5)
- OGLC (1)
- Public Domain (6)
- RL (2)
- other (39)

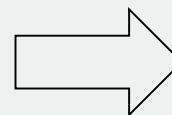
Data upload ⊕

Data upload restrictions ⊕

Automation example - DMP

DMP: *'Data is sensitive and ethical issues exist'*

DS Exercise DMP				
ID (HTTP-ORCID)	Created	Modified	Language	Description
https://doi.org/10.5281/zenodo.2644751	22.04.2019	28.06.2019	en	First and foremost, source code will be created in this experiment. This source code will be written in the programming language Python and will be partly contained in so-called Jupyter notebooks (special JSON files) and partly in python files (standard text files with the .py extension). Additionally, some documentation will be produced and saved in both a Microsoft Word file and a markdown file. The output of the source code will be stored partly in json files consisting of text columns (see ID 2 the result of classification) and partly directly in said Jupyter notebooks. Additionally, images will be produced showing a visual representation of the decision tree used for classification.
Data Officer		Name: Maria Leitinger Email: maria.leitinger@tuwien.ac.at ID: https://orcid.org/0009-0002-4537-6648 (HTTP-ORCID)		
I Data Characteristics				
I.1	Description of the data	The data in this project will be created in collaboration with Bernhard Fuhr, bernhard.fuhr@tuwien.ac.at (ID: 1 (custom)) Paul Lang - Lichtenhan, paul.lang@student.tuwien.ac.at (ID: 2 (system)) One dataset will be created. This dataset is titled "All data for project" and described as follows: First and foremost, source code will be created in this experiment. This source code will be written in the programming language Python and will be partly contained in so-called Jupyter notebooks (special JSON files) and partly in python files (standard text files with the .py extension). Additionally, some documentation will be produced and stored in both a Microsoft Word file and a markdown file. The dataset's resource type genre is "Dataset" and its language is specified as "en". Keywords for the dataset are: Breast Cancer, ArteryArtery The data will be created over the course of a single project. [Status: Embargoed]. This experiment aims to apply two machine learning models to solve a classification task for two different datasets. Each dataset is evaluated with each of the machine learning models, using different parameter settings and preprocessing strategies to compare the respective results and analyze across datasets and/or machine learning models. The project starts at 28.06.2019 and is due to end at 31.08.2019. It will be funded by 1 (custom) with the Grant ID 1-1 (custom). The funding status is currently "granted".		
II Documentation and Metadata				
II.1	Metadata standards	Just some metadata Language: en ID: 2 (system)		
II.2	Documentation of data	One technical resource is needed for this project: * ID: 0 (system) A fridge to keep the cold brew coming.		
II.3	Data quality control	* No assistance.		
III Data availability and storage				



re3data.org

Certificates

Data access

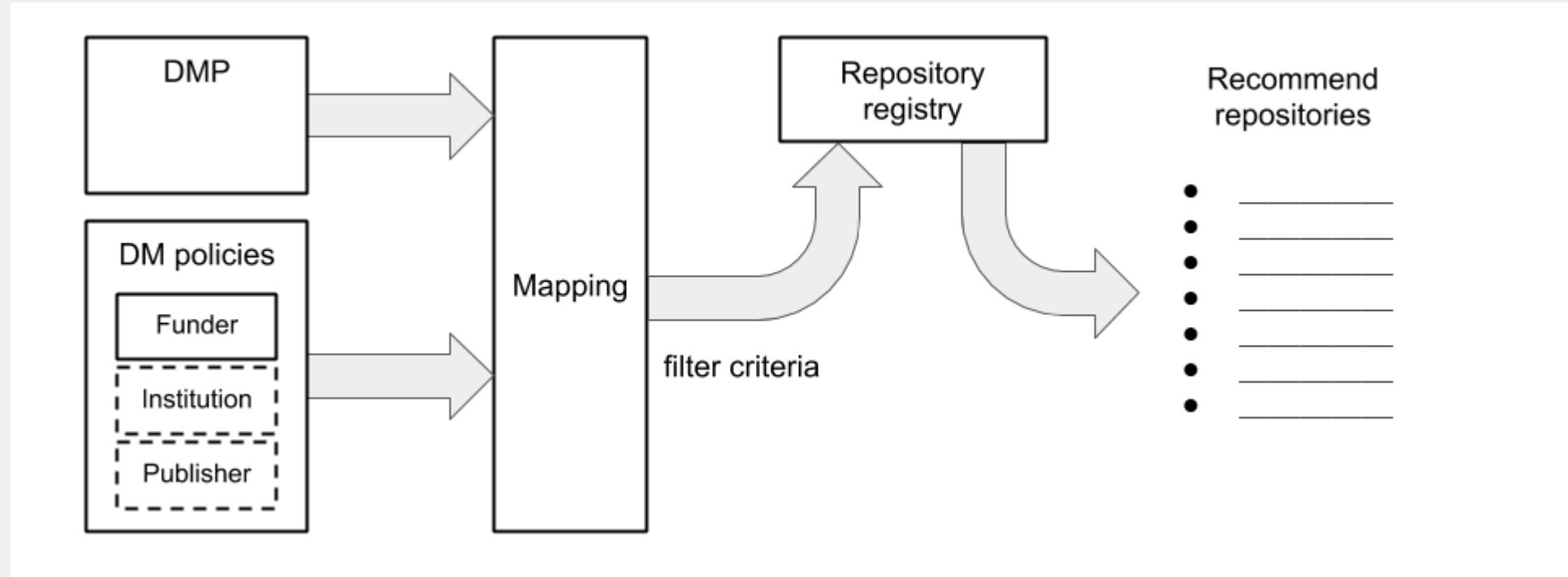
closed (161)

embargoed (62)

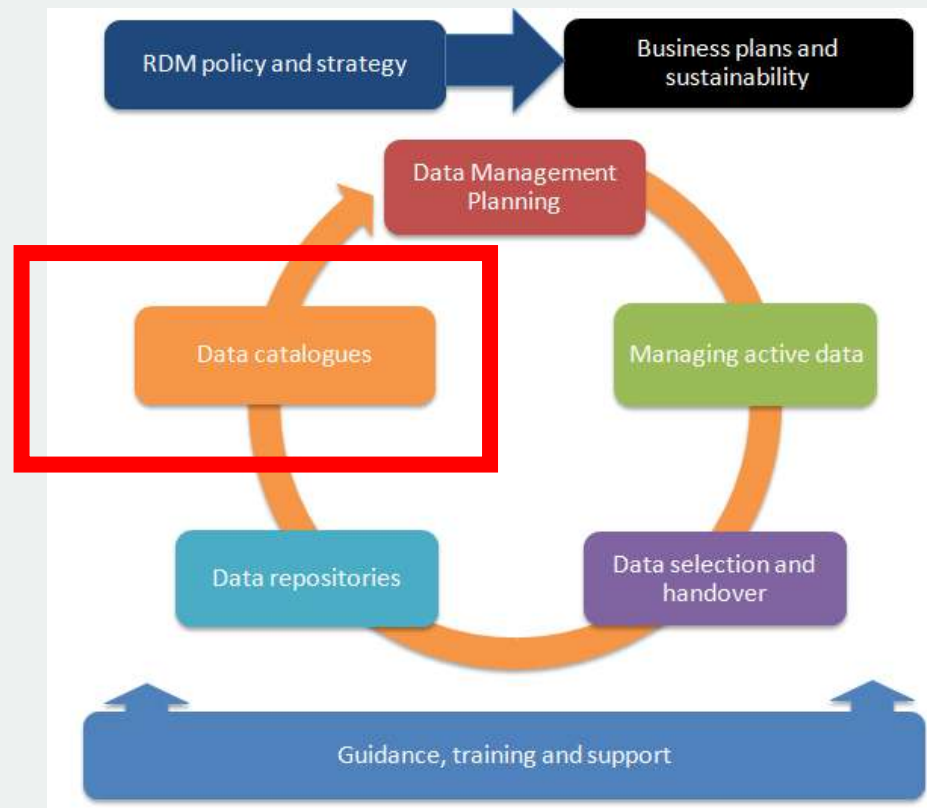
open (132)

restricted (161)

Possible solution



Simon Oblasser, Tomasz Miksa, & Asanobu Kitamoto. (2020, March). Finding a repository with the help of machine-actionable DMPs: opportunities and challenges. Zenodo. <http://doi.org/10.5281/zenodo.3701564>



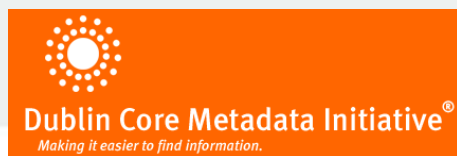
DATA CATALOGUES

Data catalogues

Topic covered in the lecture on Repositories – external visibility

Define the metadata you need to record research datasets

Expose metadata for inclusion in national catalogues or other relevant services



Summary

It is not only about technical solutions!

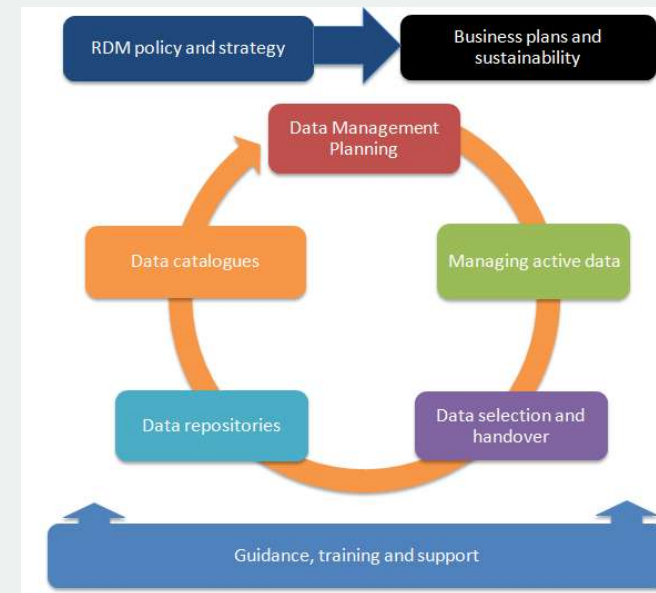
Integration is crucial

- systems
- stakeholders groups
- services

Develop a vision and a plan

Include all stakeholders

Make incremental development



REPOSITORY CERTIFICATION

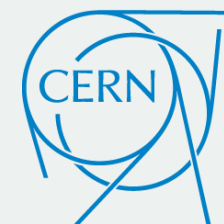
Motivation

Trust that data remains useful and meaningful into the future

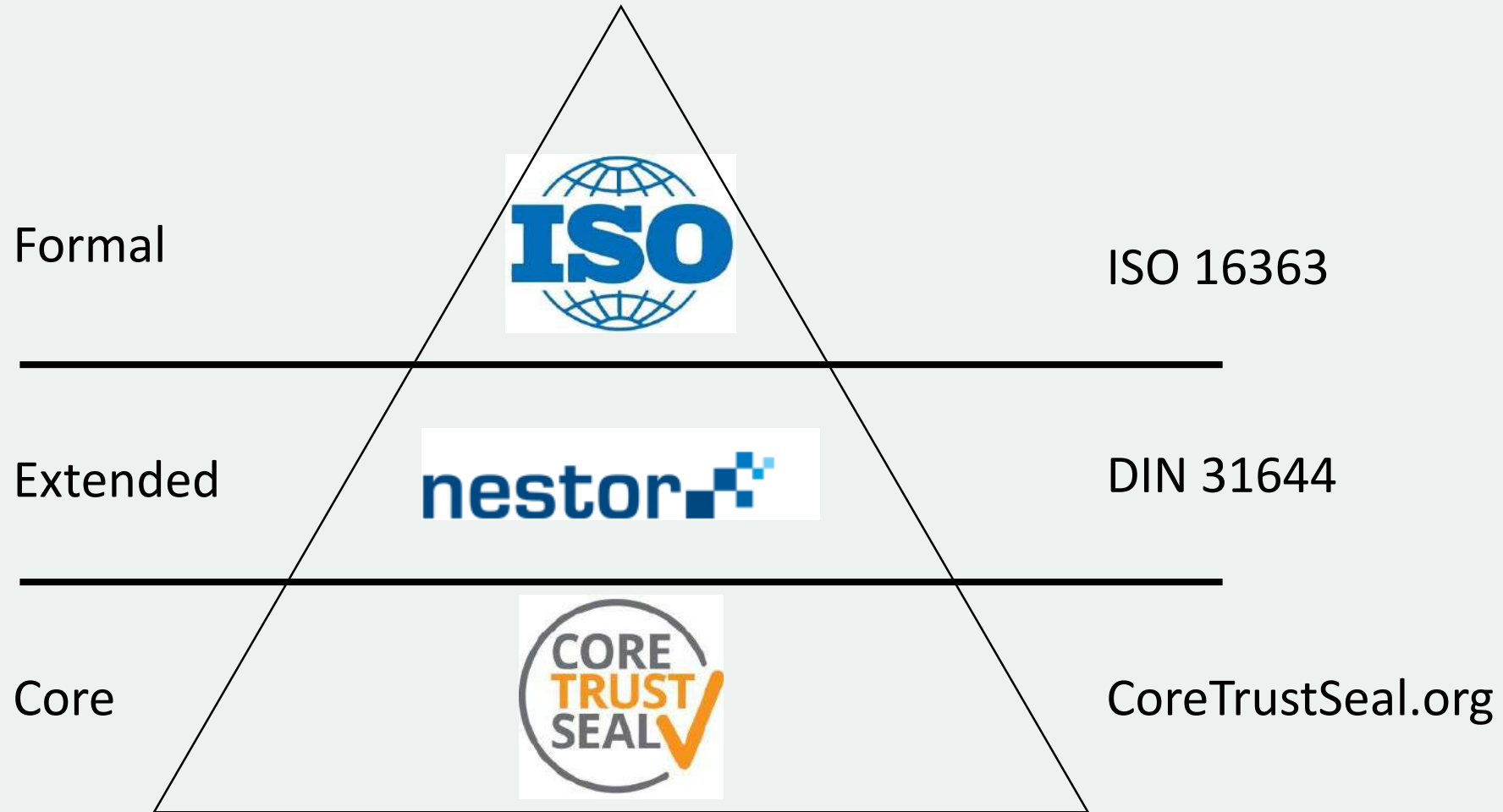
- Will the repository exist in the future?
- How is data stored?
- Is a preservation plan in place?
- etc.

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

*Are all of them trusted
repositories?*



Certification Standards



Core Trust Seal (CTS) - procedure

Self assessment based on a checklist

Guidance

- online tools
- documents and webinars

Review of the self assessment by two reviewers

Assessments publically available

Renewal every three years

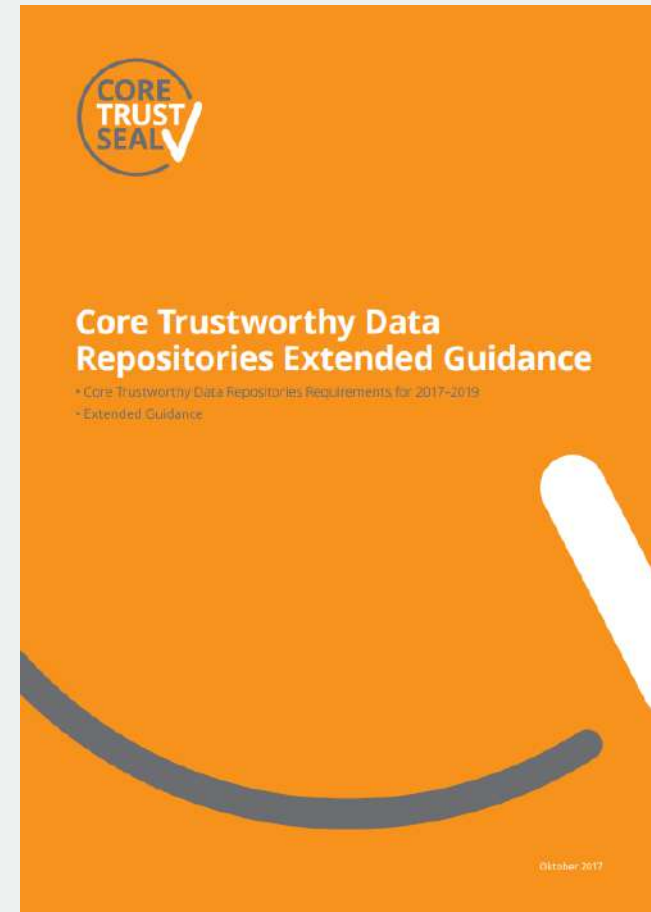


CTS Requirements

Organisational Infrastructure (6)

Digital Object Management (8)

Technology (2)



<https://www.coretrustseal.org/wp-content/uploads/2017/01/20171026-CTS-Extended-Guidance-v1.0.pdf>

Organisational Infrastructure (6)



R1. The repository has **an explicit mission** to provide access to and preserve data in its domain.


R3. The repository has a **continuity plan** to ensure ongoing access to and preservation of its holdings.

R5. The repository has **adequate funding** and sufficient numbers of **qualified staff** managed through a clear system of governance to effectively carry out the mission.

Digital Object Management (8)



R7. The repository guarantees **the integrity and authenticity** of the data. 

R10. The repository assumes responsibility for **long-term** preservation and manages this function in a planned and documented way. 

R12. Archiving takes place according to **defined workflows** from ingest  to dissemination.

R13. The repository enables users to **discover the data** and **refer to them in a persistent way** through proper citation. 

Technical (2)

R15. The repository functions on **well-supported operating system** and other core infrastructure software is using hardware and software technologies appropriate to service it provides to its Designated Community.

R16. The technical infrastructure of the repository provides for **protection** of the facility and its data, products, services, and users.

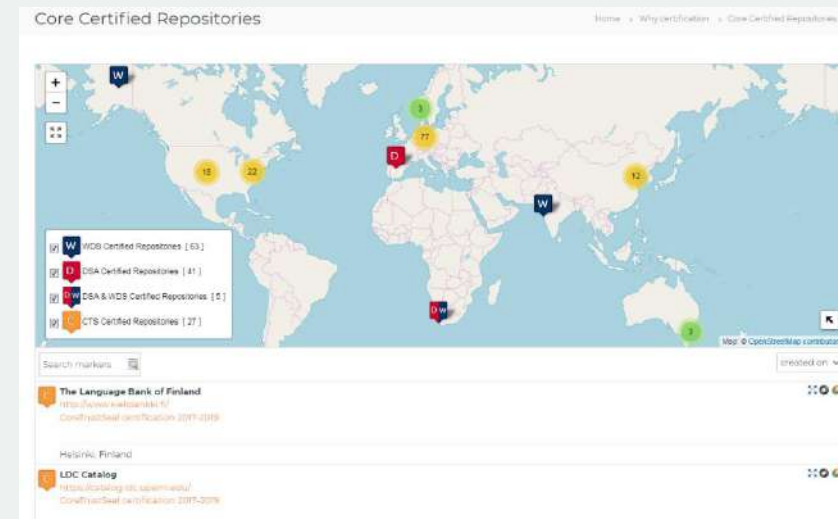
Self-assessment examples

Austrian example

- A Resource Centre for Humanities Related Research in Austria
 - <https://arche.acdh.oeaw.ac.at>
- Self assessment
 - <https://www.coretrustseal.org/wp-content/uploads/2018/03/ARCHE.pdf>

More examples

- <https://www.coretrustseal.org/why-certification/certified-repositories/>



X. Preservation plan

R10. The repository assumes responsibility for long-term preservation and manages this function in a planned and documented way.

Compliance Level:

Response

Guidance:

The repository, data depositors, and Designated Community need to understand the level of responsibility undertaken for each deposited item in the repository. The repository must have the legal rights to undertake these responsibilities. Procedures must be documented and their completion assured.

For this Requirement, responses should include evidence related to the following questions:

- Is a preservation plan in place?
- Is the 'preservation level' for each item understood? How is this defined?
- Does the contract between depositor and repository provide for all actions necessary to meet the responsibilities?
- Is the transfer of custody and responsibility handover clear to the depositor and repository?
- Does the repository have the rights to copy, transform, and store the items, as well as provide access to them?
- Are actions relevant to preservation specified in documentation, including custody transfer, submission information standards, and archival information standards?
- Are there measures to ensure these actions are taken?

Self-assessment example

10. Preservation plan

Applicant Entry

Statement of Compliance:

3. In progress: We are in the implementation phase.

Self-assessment statement:

As its primary preservation strategy, ARCHE performs migration of formats as opposed to providing software emulation. It aims to establish a high level of transparency with its depositors and users. Thus the deposition agreement and other relevant informative sections of our website highlight our responsibilities and our rights to copy, transform, store and provide access to the deposited items. All the actions relevant to preservation are specified in our documentation.

(this is just an excerpt of the full answer)

<https://www.coretrustseal.org/wp-content/uploads/2018/03/ARCHE.pdf>

Advice from Uni Graz



one team lead

collaborative document to work jointly

put procedures and workflows in writing

enhance publicly available information

use the glossary and extended guidance

- 'we' – 'the repository' – 'the system'

read successful applications on website

allow sufficient time for the process of internal coordination and data gathering

technical procedures do not pose problems/attract inquiries from reviewers, it is organizational issues that do

Based on presentation by Elisabeth Steiner at the Certification Workshop on fair-aligned repositories in Austria, 14.11.2019, Vienna

NESTOR

DIN 31644

- Deutsches Institut für Normung
- 'Criteria for trustworthy digital archives'

Derives from Trustworthy Repositories Audit & Certification (TRAC)

- document describing the metrics of an [OAIS](#)-compliant digital repository
- TRAC is discontinued, replace by the ISO 16363
- translated to German

Uses OAIS terminology

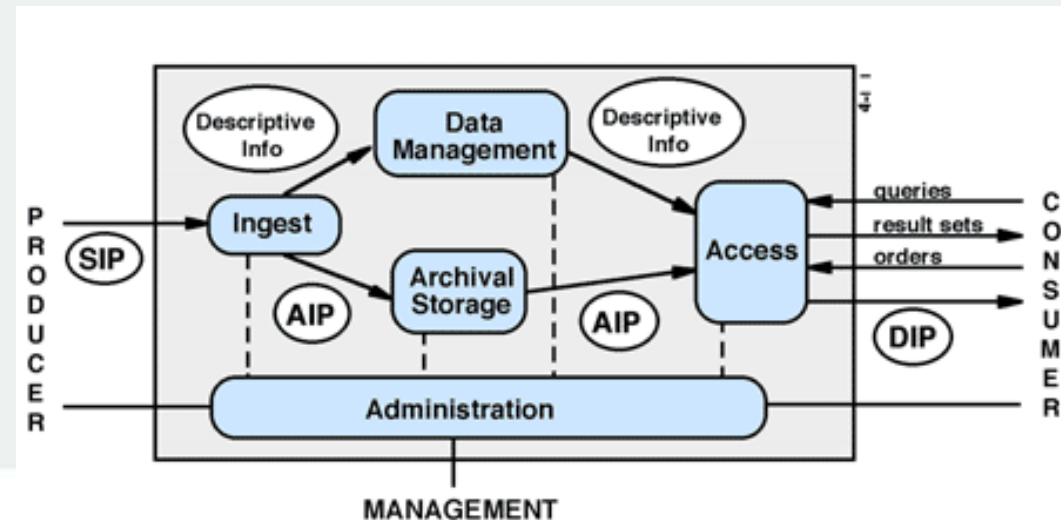
The logo for nestor, featuring the word "nestor" in a bold, blue, sans-serif font. To the right of the text is a graphic element consisting of several blue squares of varying sizes, arranged in a pattern that suggests a staircase or a digital grid.

Open Archival Information System (OAIS)

Discussed in a separate lecture

OAIS

- conceptual framework for an archival system dedicated to preserving and maintaining access to digital information over the long term
- defines concepts, terminology
- not to be instantiated!



NESTOR requirements overview

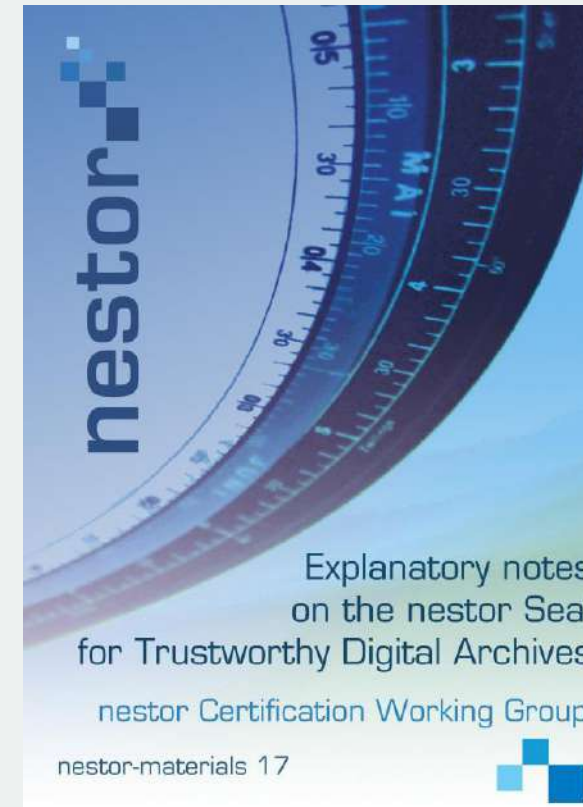
34 criteria

- 12 obligatory
- other can be excluded when justified

auxiliary questions

self-assessment

2 external reviewers



<https://d-nb.info/1047613859/34>

C27 Identification

A digital archive should use internal identifiers to manage the information objects and their representations and, where applicable, their parts and relationships (part/totality, different variants, versions etc.), especially to ensure unique assignment of the content data to the metadata.

The use of externally visible, standardised persistent identifiers ensures reliable tracing of the information objects and their representations, and consequently also access.

To what extent must the criterion be met?

An average of 7 points must be achieved in the assessment of the applicable criteria C13 - C34.

Explanation: The information objects, representations and their parts are permanently linked to each other. These links can only be preserved through the use of persistent identifiers. The identifiers should not change over the course of time (i.e. be permanent) and should be created using uniform specifications. They should be recognisable to external users, producers and others. By entering the identifier, external users should be able to find and use the required object. Possible specific requirements for identifiers are described e.g. in DIN 13646 "Requirements for the long-term handling of persistent identifiers".

Questions

- Which identifiers does the digital archive use?
- Which procedure has been used to give unique identifiers to all information objects, representations and their parts, and to all content and metadata?
- How is the identifier-based assignment conducted?
- How is the permanence of the identifiers ensured?
- How are the identifiers made available to external users?

Documents: Specification of the internal and external identifiers

ISO 16363

Uses OAIS terminology

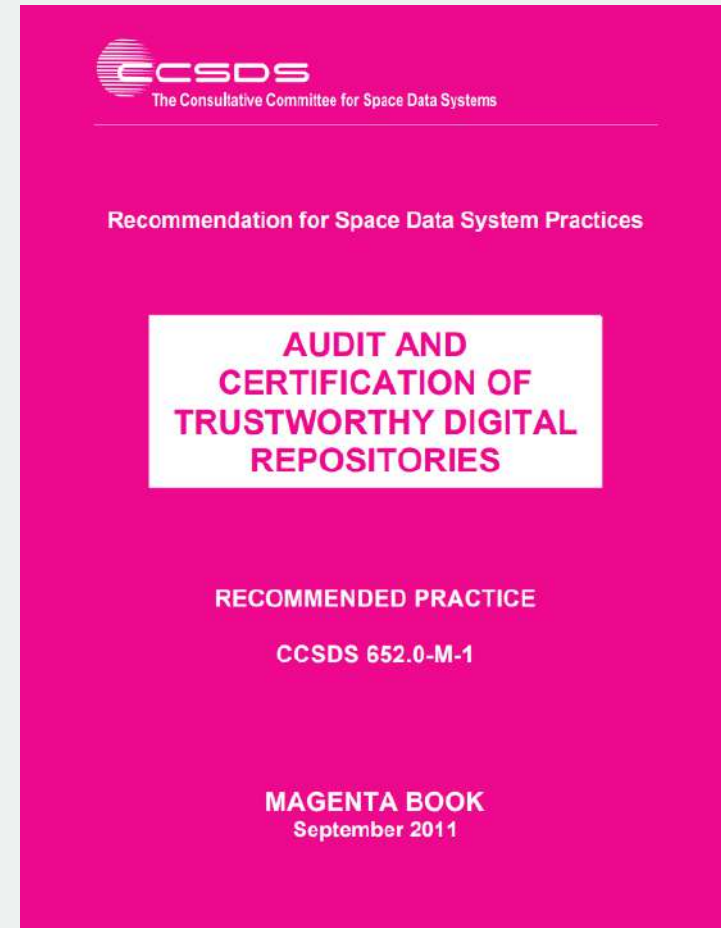
Based on TRAC

Over 100 metrics

Detailed instructions

Full external auditing process

ONE repository
certified since 2012!



<https://public.ccsds.org/Pubs/652x0m1.pdf>

Structure

Organizational infrastructure

- Governance and organizational viability
- Financial sustainability
- ...

Digital object management

- Ingest: acquisition of content
- AIP Preservation
- ...

Infrastructure and security risk management

- Security risk management
-

Metrics

Metrics and their structure:

- Statement of requirement
- Supporting text
- Examples of Ways the Repository can Demonstrate it is Meeting this Requirement
- Discussion

4.3 PRESERVATION PLANNING

4.3.1 The repository shall have documented preservation strategies relevant to its holdings.

Supporting Text

This is necessary in order that it is clear how the repository plans to ensure the information will remain available and usable for future generations and to provide a means to check and validate the preservation work of the repository.

Examples of Ways the Repository Can Demonstrate It Is Meeting This Requirement

Documentation identifying each preservation risk identified and the strategy for dealing with that risk.

Discussion

These documented preservation strategies will describe how the repository will act upon identified risks, as part of the preservation strategic plan. These preservation strategies and the preservation strategic plan will typically address the degradation of storage media, the obsolescence of media drives, and the obsolescence or inadequacy of Representation Information (including formats) as the knowledge base of the Designated Community changes, and safeguards against accidental or intentional digital corruption. For example, if

Metrics – self assessment

4.3 PRESERVATION PLANNING					
	Metric	Supporting Text	Examples of Documents the Repository can use to demonstrate it is Meeting this Requirement:	Brief description of evidence (add rows if necessary to list all relevant documents for a metric) Use short titles for documents. Provide detailed	Explanation of how the repository addresses this metric
4.3.1	THE REPOSITORY SHALL HAVE DOCUMENTED PRESERVATION STRATEGIES RELEVANT TO ITS HOLDINGS.	This is necessary in order that it is clear how the repository plans to ensure the information will remain available and usable for future generations and to provide a means to check and validate the preservation work of the repository.	Documentation identifying each preservation risk identified and the strategy for dealing with that risk.		
4.3.2	THE REPOSITORY SHALL HAVE MECHANISMS IN PLACE FOR MONITORING ITS PRESERVATION ENVIRONMENT.	This is necessary so that the repository can react to changes and thereby ensure that the preserved information remains understandable and usable by the Designated Community.	Surveys of the Designated Community of the repository.		
4.3.2.1	<i>The repository shall have mechanisms in place for monitoring and notification when Representation Information is inadequate for the Designated Community to understand the data holdings.</i>	This is necessary in order to ensure that the preserved information remains understandable and usable by the Designated Community.	Subscription to a Representation Information registry service; subscription to a technology watch service, surveys amongst its Designated Community members, relevant working processes to deal with this information.		

http://www.iso16363.org/?smd_process_download=1&download_id=30

ISO Process for Audits

Preparatory work by repository

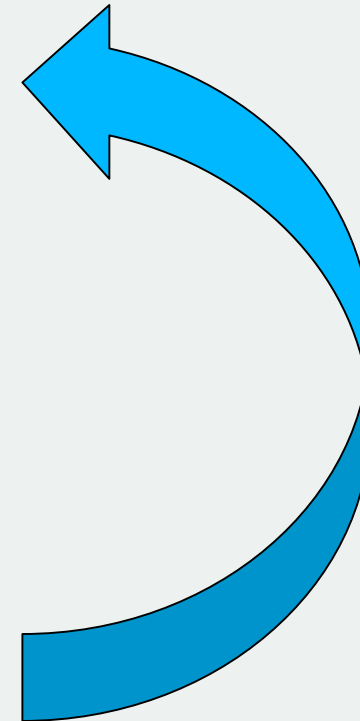
First audit and resulting certification

- Identifies improvements needed
- Repository prepares improvement plan

Repository implements improvement plan

Surveillance audit after a period

Re-certification



Why certify a repository (using any scheme)?

Builds stakeholder confidence in the repository

- researchers', funders', publishers', etc.

Improves communication within the repository

- roles and responsibilities (implicit knowledge becomes explicit)
- processes

Ensures transparency

- most assessments are public

Enables comparison of repositories

Attracts funding

Summary

You should know

What are the lifecycle models and how to use them?

What are the components of an RDM service infrastructure?

What is the scope of policies and how they drive DM activities and obligations?

How do develop support services and who are data stewards?

What to consider when implementing DMPs in an institution?

What standards exist and why certification matters?

What are the certification criteria?

Acknowledgments

Jones, S., Pryor, G. & Whyte, A. (2013). 'How to Develop Research Data Management Services - a guide for HEIs'. DCC How-to Guides. Edinburgh: Digital Curation Centre. Available online: <http://www.dcc.ac.uk/resources/how-guides>

Andreas Rauber. (2017, November). Research Outputs Management: what services are involved? File (Version 1.00). Zenodo. <http://doi.org/10.5281/zenodo.1063548LEARN>

Teperek, Marta. (2017, November). Data Stewardship - addressing disciplinary data management needs at TU Delft. Zenodo. <http://doi.org/10.5281/zenodo.1064794>

Budroni, Paolo. (2017, May). The LEARN Project Using the LEARN RDM Policy & Guidance. Zenodo. <http://doi.org/10.5281/zenodo.579993>



TECHNISCHE
UNIVERSITÄT
WIEN



data and legal aspects

Verena Dolovai – Compliance Officer Research



The word 'AGENDA' is written in large, white, bold, sans-serif capital letters. It is centered and surrounded by a cluster of overlapping, semi-transparent blue squares of various sizes. To the right of the word, there is a small, dark blue square containing the white text 'TU' above 'WIEN', which is the logo of TU Wien.

- **types of data**
- **ownership of data**
- **inventor, author, creator**
- **fields of law**
- **rights in research data**
- **copyright in research data?**
- **rights holders in research data**
- **data protection law & research project**
- **examples**

- instrument measurements
- experimental observations
- still images, video and audio
- text documents, spreadsheets, databases
- quantitative data (e.g. survey data)
- survey results & interview transcripts
- simulation data, models & software
- slides, artefacts
- specimens, samples, questionnaires



- owner

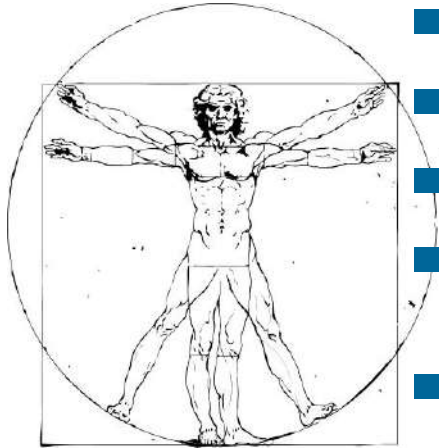
- a term of civil law
- not exactly suitable for data ownership
- however, „owner“ has become familiar in terms of data management



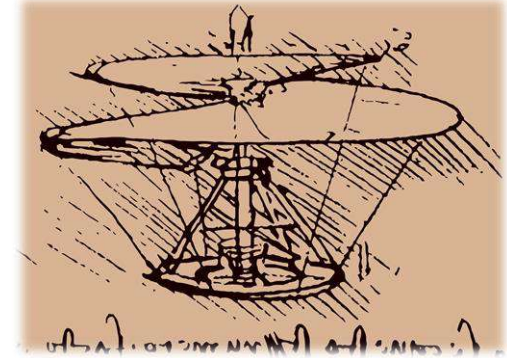
- rather, legally speaking:

- author (copyright law)
- creator (if not protected by copyright law)
- inventor (patent law)
- legal successor





- **inventor → patent law**
- patent law: protects inventor, idea, technology
- -> exclusive right, registration, fees
- 20 years protection



- **author → copyright law**
- copyright law: protects author, work, intellectual creation (works of literature, science, art, databases, computer programmes)
- -> no registration, no fees
- 70 years protection

- **creator (person who creates data(sets) → licenses**
- licenses: protect creator





- copyright law (databanks, images)
- data protection law
- data security law
- competition law (protection of know-how)
- contract law
- labour law (towards employees of an institution)



data protection law in research

rights in research data

GDPR*: personal data
& informed consent

general legal aspects of
research data

*as well as national data
protection laws: Austrian
Data Protection Act (DSG)
and Research
Organisation Act (FOG)

licenses, open source
software, IP, storage,
erasure, publication, user
rights in data, user
agreements,
confidentiality**, liability etc)

**data secrecy and protection
of confidential information



rights in research data

data*

copyright and neighbouring
rights
data protection law

software

copyright

copyright relevance:
**qualitative research
data**; data in social
sciences and humanities

database

right in the database (work)

data medium

ownership

*as a rule, data in a
literal sense are
hardly subject to
copyright; however,
term „data“ must be
defined in each
individual case!

copyright in research data?

essays
articles
papers
monographs
graphics
maps
tables with individual graphic design
databases
research software
photographs
video recordings



possible copyright:
result of individual
creation

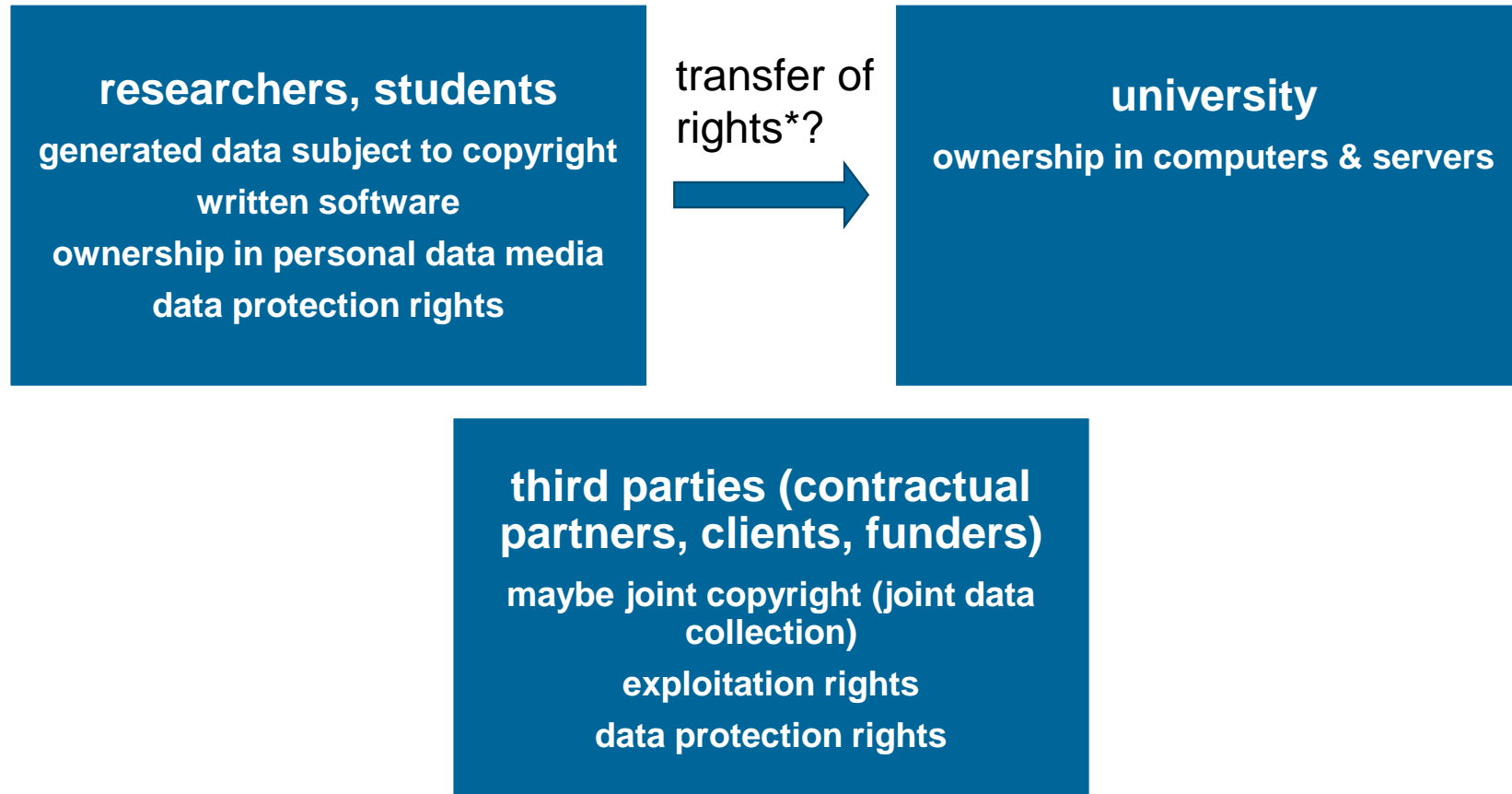


measurement data
metadata
results of a software simulation



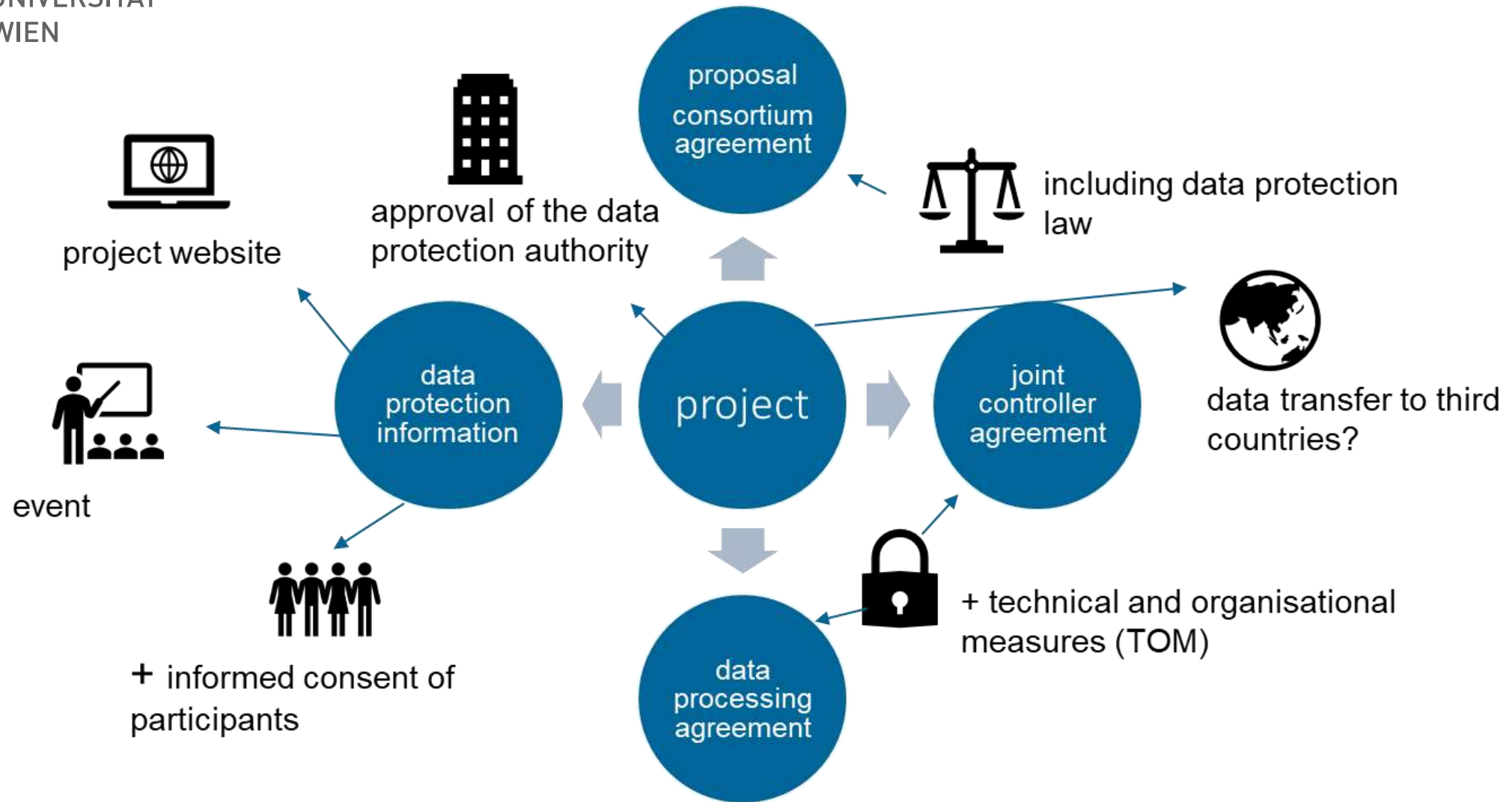
usually no copyright:
no result of individual
creation

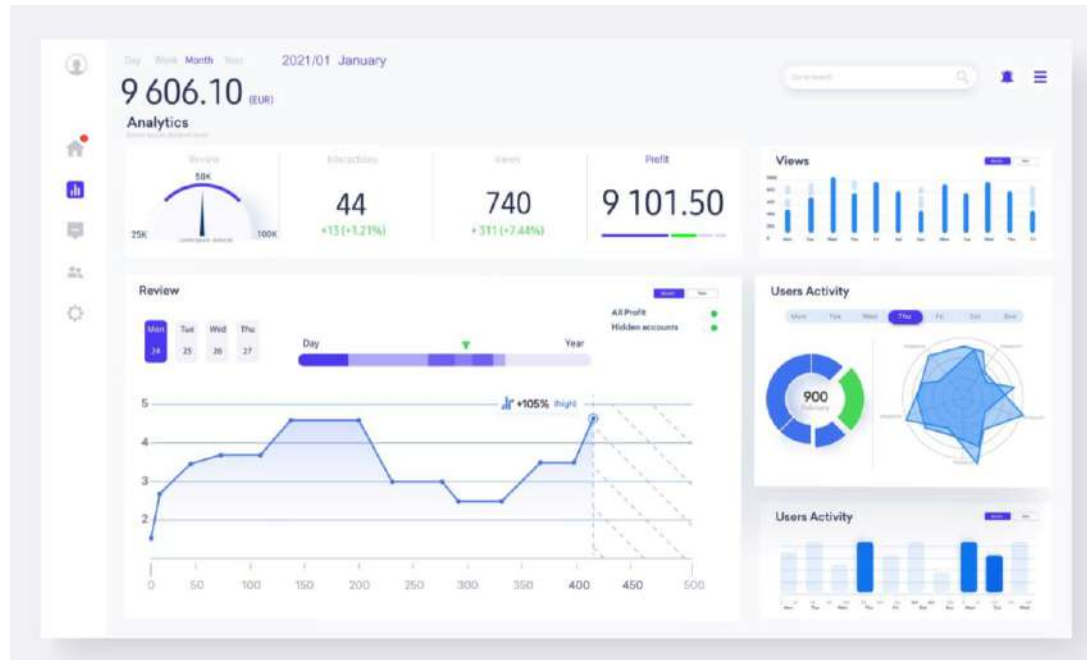
**case by case
assessment!**



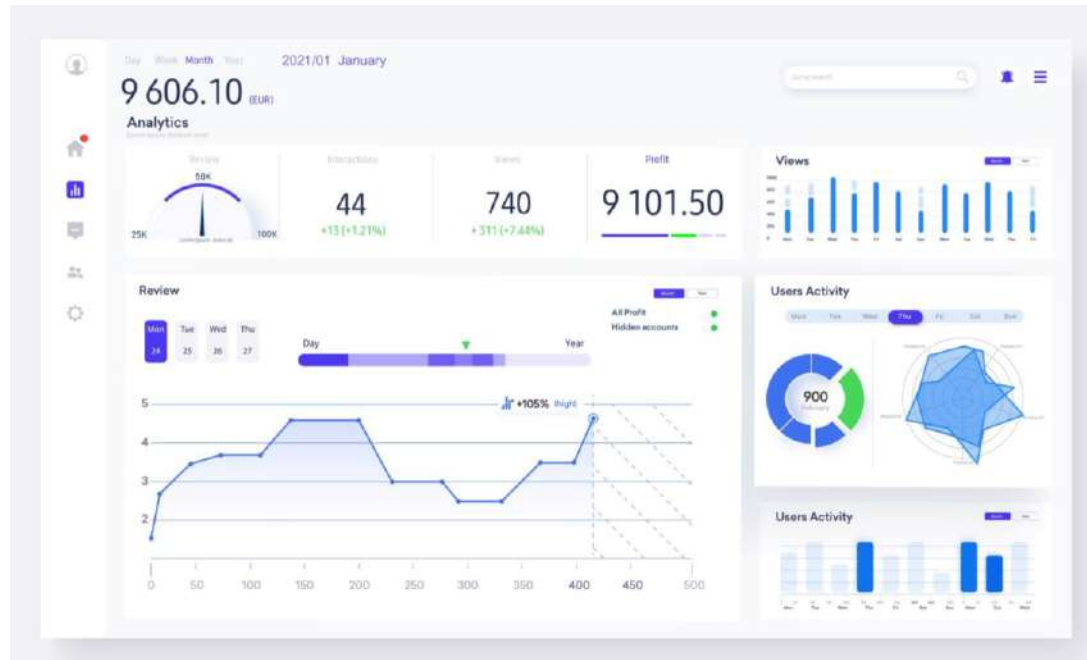
*employment contract, exploitation rights, and user rights

data protection law & research project





I want to use data from a website



- I want to use data from a website
- check for license or other legal requirements / rules
 - not available?
 - seek legal advice
 - include reference



I'm leaving my organisation
What to do with „my data“?



I'm leaving my organisation What to do with „my data“?

- check employment contract
- check project contract / funding agreement
- seek legal advice & agreement with employer



I am employee of xx and I want to license „my data“. Am I free to choose any license?



I am employee of xx and I want to license „my data“. Am I free to choose any license?

- check for internal guidelines (types of licenses, powers of attorney, regulations)
- check employment contract
- check project contract / funding agreement



**Someone used my data illegally
or without my permission...
what should I do?**



**Someone used my data illegally
or without my permission...
what should I do?**

- **seek legal advice**

Verena Dolovai
Compliance Officer Research
Technische Universität Wien
Favoritenstraße 16/DG
1040 Wien
Telefon: +43 1 58801 406635
verena.dolovai@tuwien.ac.at

