

## Love your data!

Data Management Basics for study and research

12.08.2020

Timo Gnadt gnadt@sub.uni-goettingen.de



eResearch



## Survey: Who are you? (1)



- A: I am a bachelor student
- B: I am a master student
- C: I am a doctoral/PhD student or scientific assistant
- D: I am a PostDoc
- E: None of the above



## Survey: Who are you? (2)

#### I am from the ...

- A: natural/structural sciences
- B: arts and humanities
- C: social sciences
- D: Jurisprudence
- E: None of the above





## Survey: Who are you? (3)

I am studying / doing a PhD / working

- A: in Göttingen
- B: not in Göttingen





## Structure

- Introduction to Research Data Management
- Storing data
- Organizing data
- Describing data
- Sharing data
- Further links and information

Comments, questions and suggestions for this training can be entered here:

https://pad.gwdg.de/LoveYourData\_20200812?both



## Introduction to Research Data Management

## eResearch

## What is 'data'?



## What is 'data'?

"A re-interpretable representation of information in a formalised way, suitable for communication, interpretation or processing."

Data are representations of observations, objects or other units used as evidence of phenomena for the purpose of research or scholarship.

#### (Christine Borgmann, UCLA, 2014)

translated into Digital Curation Centre



A second se





## What is research data?

#### All representations of information that you use in your research:

**Types:** Statistics, interviews, simulations, measurement data from experiments, observation data from instruments, text with semantic annotations, 3D scans, model drawings, numerical representations, ...

**Shapes:** Video, audio, images, spreadsheets, e-mails, paper documents, binary data, software, text files, lab notebooks, ...



Graph from: Ouborg NJ, Pertoldi C, Loeschcke V, Bijlsma R, Hedrick PW (2010). <u>Conservation</u> <u>genetics in transition to conservation genomics</u>. Trends Genet 26: 177-187.

## What is *your* research data? or: What will *your* research data be?

#### Answers from the participants:

- genomic data
- metric data derived from GIS grids
- · Archival materials
- · electrophysilogical recordings
- · Video data
- neuroimaging (MRI)
- · observation from field testing
- · Genomic data
- · observational data from biodiversity surveys
- material and geological data
- · Observations from greenhouse plants
- Fluorescence, electrophysiological data
- · Oral History data and Archival materials
- Manuscripts
- field survey

Alliance



### Types of research data

Туре	Characteristics	Example
Observations	Data is collected in real time Mostly irreplaceable	Sensor data Survey data
Experiments	Mostly created in the laboratory Reproducible but expensive	Gene sequences Chromatogram
Simulations	Generated from test models Model and metadata more important than output	Climate models Economic models
Derived data	Derived or compiled from other data, reproducible	Text Mining 3D models
References	Collection of smaller data sets Mostly published	Gene sequence database Primary text sources
Digital copies	Digital version of an analog object, reproducible as long as the original exists	Manuscripts

### Research data - a valuable investment



Source: European Space Agency: Rosetta and Philae at comet, on flickr.

#### Rosetta Mission, 2004-16

**Duration:** 

- >10 years preparation
- 10 years from the beginning to the data

eResearch

costs:

• over € 1.000.000.000

#### Result:

- some cool photos
- many data
- a radically new theory about the origin of the universe?

Alliance

### Research data - a resource in great demand



#### Data with reference to COVID-19

- medical, molecular biological, epidemiological...
- mathematical, sociological, geographical, psychological, ...

Alliance

eResearch



### Research cycle





#### Answers from the participants:

- · propper documentation and sustainable storage
- · data collection and classification
- sorting out the methods for saving data in the way to then easily find it. Also systematise the data according to the content or technical characteristic..
- · saving and making data accessible to others
- · finding the digital space to store it

Alliance

eResearch



## eResearch

# Research data policy of the Georg-August-University of Göttingen

- Officially published on August 28, 2014
- One of the first German universities with such a guideline
- Topics covered:
  - Research data, research data management and its purposes
  - Data management plans
  - Support, training and service delivery
  - Storage solutions
  - Ethical and legal standards
  - Open Access
- eResearch Alliance: Support and consultation for the implementation of the policy for the Göttingen Campus

Source: http://www.uni-goettingen.de/en/488918.html



The Goop-Jugat-University Coefficient is committed billigently preserve results of scholarship, to produce never results through research, and to make results accessible and results for scheduria and the wider societies, never after that regenerations. The management protection, preservation and sustainable provision of research data must therefore be carried out in accordance with recognized standards, meet high expectations and full legal and ethical obligations. The University accounsel standards are that the implementation of this judities will depend on the settings and requirements of each subject area.

University promotes and supports open access to research data.

Research data are those data collected, observed, simulated, derived, or generated during the course of research.

inagement of research data includes their planning, collection, processing, and preservation. It ensures the access to, and the reuse, reproducibility, and quality assurance all research data underplinning research results.

I. Research data management is generally the responsibility of the person leading a project and the researcher who is acting in an individual capacity. A particular responsibility is the adherence to good practices of research as well as standards in their subject area.

5. Research projects with research data require a data management plan that includes but is not restricted to the topics of access rights to research data and necessary precautions for handling them.

 The University provides support and advice for research data management in the preparatory stages of research projects, during their conduct and after their completion, and provides appropriate training.

 The University implements and maintains essential services for research data infrastructure that ensures adequate storage and technical availability of digital research data Specific requirements have to be aligned among all stakeholders and may involve additional funding.

8. Storage and archiving of digital research data is carried out within the technological and informational infrastructure of the University or in acknowledged external or internal subject repositories.

3 The University and its researchers of their research data management to given conditions of ethics, data protection, intellectual property privacy and disclosure. This leaves regulations unbuched that relate to an assessment of research data according to the German employee invention act and specific contractual agreements.



- possibility of re-analysis
- · collect data for a larger context than own data

- Increase your efficiency and the value of your work / research
  - Prevent data loss
  - Avoid unnecessary work
  - Improve data quality

#### Good scientific practice

- Reproducibility, accountability and compliance
- "Primary data as a basis for publications must be stored securely in the institution where they were created for ten years in a permanent form". (DFG, Proposals to safeguard good scientific practice, 1998)
- Request from funding organisations (e.g. EU Horizon 2020)

#### • Data exchange with colleagues

- Research can be *very* expensive, and the only result of long research projects and trips is often data.
- Compared to the costs of data collection, the costs of data management are low.
- Productive data exchange is simply a matter of efficiency.

Alliance

eResearch



eResearch / Alliance









#### Sources:

Jeffrey Brainard et al., **Rethinking retractions**, *Science* 26 Oct 2018: Vol. 362, Issue 6413, pp. 390-393 DOI: 10.1126/science.362.6413.390 GRAPHIC: J. YOU/SCIENCE; DATA: RETRACTION WATCH

16.06.2020



- Improve your research
- Good scientific practice
- Data exchange with colleagues
- Data Publication
  - Required by increasing number of journals
  - Avoid retractions
  - Get cited for your data!
- Enabling new types of research
  - Feedback loops between empirical and modelling approaches
  - Initiation of research questions in completely different areas



## Publications are arguments of authors, and data are the evidence to support these arguments.

(Christine Borgmann, UCLA, 2014)

# The deeper meaning of Research Data Management





Alliance



## Levels of data curation





## Storing data



### Why do I want to store my data?

	Purpose	Solution
?	short-term storage, backup copy	<ul><li>institutional backup</li><li>individual backup</li></ul>
	long-term storage, archiving	<ul> <li>Data archive,</li> <li>institutional archive solution or</li> <li>individual archiving</li> </ul>
	making data available to others	<ul><li>Data repository:</li><li>institutional,</li><li>generic or</li><li>subject-specific</li></ul>



## Yes, we store – what for?

	Backup	Archiving	Sharing
Purpose	Ability to recover data in case of data loss or error propagation	Enabling validation through persistent storage of data used in publications	Enable review, citation and reuse of data sets (data exchange)
Characteristics	Duplication of current work data & interim results	Archive format (e.g. zip) with all associated & relevant data / files (ideally including metadata)	Format defined by the repository; discipline- specific metadata standards
Regularity	Regularly during the work phase or project duration	Once for each relevant data set, usually at the end or after the work phase	Once for each selected data set, either during or after the work phase
Expenses	Varying, e.g.: set up once, check regularly	Establish predefined procedure with data archive (e.g. data/computing center)	Documented process, sometimes accompanied by repository support



## Survey: Backup 1 Do you back up your research data? If so, how?

- A: No, never.
- B: Yes, on a USB stick.
- C: Yes, on an external hard disk or CD/DVD.
- D: Yes, in the cloud.
- E: Yes, through a service of my institute/university.





Alliance

## Survey: Backup 2 How often do you back up your research data?

- A: Not at all
- B: 1-2 times a year
- C: 3-12 times a year
- D: Several times a month
- E: (Almost) daily

#### Answers from the participants:



eResearch



## Survey: Backup 3

Have you ever tried to recover a deleted file?

- A: No, never.
- B: Yes, but not successfully
- C: Yes, successfully but with high effort
- D: Yes, that's no problem for me



## Survey: Backup 4

Can you revert to an earlier version of a file?

#### A: No, I don't know how.

- B: Yes, I think so.
- C: Yes, I've done that before.
- D: Yes, that's no problem for me.





## Survey: Backup 5

## Who is responsible for backup and storage services at your institute?

- A: I don't know.
- B: I think I could find out.
- C: I know who I can ask about that.
- D: I have already had contact with them.





## Why backup?

#### Notice at the bus stop Jüdenstraße/Göttingen


#### Why backup?

Source: Gino on flickr

Source: steviep187 on flickr 😳 🤨 💬

Source: <u>Kiell Eson on flickr</u> Source



#### Why backup?



**Source:** University of Southampton, Department of Electronics and Computer Science, 2005

Source: Frankfurt University of Applied Sciences,

March 2020 © Frankfurt Fire Department



#### Causes for data loss



Answers from the participants:

- virus attack
- errorness deleting
- messy storing
- laptop crash

### Causes for data loss

- Malware / theft / destruction
- Software failures
  - Program errors / bugs / software updates
  - Features

(e.g.: overwriting in Dropbox during synchronization)

- Hardware failures
  - Bad design / cheap parts / defects
  - Age
  - Dropped laptops / HDDs
  - Liquids (water, coffee, coke)
  - Lightning strikes / electrical impulses
- Human error
  - Accidental deletion
  - Lack of knowledge

16.06.2020

Further reading : disasters and tales of data loss, statistics on how data gets lost



eResearc

Source: a man working at home while eating breakfast by Socialeurope via flickr



Source: Kroll Ontrack, 2007, Robin Harris

Alliance

## **Backup Principles**

- Create multiple backups
- Expect human errors (keep older versions)
- Do not use backup drives for file sharing
- Store backups physically separate from your PC/laptop
- Check your backups regularly
- Practice the worst case and do a full system recovery
- Discuss the topic with friends to learn about their best practices
- Include your mobile devices in your planning

3 copies

- 2 different media
  - 1 remote storage location

eResearc

### Backup software



Operating system	perating Integrated vstem Backup-SW		Comment					
Windows 7	File Recovery		<ul><li>Requires configuration to copy not only local libraries</li><li>Can create a bootable image</li></ul>					
Windows 8/10	File History ("File version history")		<ul> <li>Saves only local libraries</li> <li>Can be configured for individual libraries and folder exclusion</li> <li>Cannot create a bootable image</li> </ul>					
Mac OS	Time Machine		<ul> <li>Saves everything except excluded folders</li> <li>Can use encryption</li> <li>Can be used to restore a nonbooting Mac</li> </ul>					
Ubuntu	File Re	ecovery	<ul><li>Uses encryption and compression</li><li>Can use cloud storage</li></ul>					
Operating syst	em	Free othe	Free other backup software					
Windows		Personal Backup, PureSync, Paragon Backup&Recovery, Robocopy,						
Mac OS		Carbon Copy Cloner, SuperDuper,						
Ubuntu		Rsync, Back in Time						

16.06.2020



### GWDG storage services

Name	Backup	Shared Usage	Comment
Fileservice / Active Directory	Yes	Possible	<ul><li>Network drives, e.g. P: , possibly more</li><li>are saved automatically</li></ul>
IBM Tivoli Storage Manager (TSM)	Yes	No	Offer for institutes to centrally secure all local workstations
CrashPlanProE	Yes	No	<ul> <li>Individual backup solution</li> <li>GWDG license: €26,- per year</li> </ul>
ownCloud	Yes	Yes	Free storage space: 50 GB
GRO.data (Dataverse)	Partially	Yes	Primarily for data exchange and data storage for later publication
HSM	No	No	For archiving data from completed projects
GitLab	No	Yes	Versioning; not for large data volumes



# Organizing data



## Why organize?



Source: austinevan on flickr

Organize your files so that you and your colleagues can find and access things when you need them.



## Why organize?



## File naming conventions

To stay organized, you should define the following:

- A self-describing folder structure or a tagging scheme
- What information should be included in file names?
- How should file names be structured?
- · How should files be referenced and how should they be exchanged?

... especially when working in a team!

Self-explanatory file name:

20200812\_Presentation\_eRA\_v42.pptx

vs. short file name:

eRA\_final.pptx

16.06.2020

Avoid special characters

eResearc

Alliance

### File naming conventions

To stay organized, you should define the following:

- A self-describing folder structure or a tagging s
- What information should be included in file nak.
- How should file names be structured?
- How should files be referenced and AND STICK

... especially when working in a team!

Self-explanatory file name:

20200812\_Presentation\_eRA\_v42.pptx





16.06.2020

#### Avoid special characters

eResearc

Use what works for

exchange





## Versioning

20200812\_Presentation\_eRA\_V13.pptx 20200812\_Presentation\_eRA\_V13final.pptx 20200812\_Presentation\_eRA\_V13new-final.pptx 20200812\_Presentation\_eRA\_V13final-finalv1.pptx 20200812\_Presentation\_eRA\_revised\_v01a.pptx

Best practice:

- Use consecutive version numbers and initials of the author
  - No unreliable descriptors in file names like "final"
  - Rather use folders for marking/sorting various purposes and for structuring
- When you work with others, agree on a common naming system
- AND/OR: Use a file versioning system
  - e.g. gitlab or ownCloud

## Version control with git

- Mainly used in software development
  - many functionalities for the support of development processes
- Also usable for versioning of documents

#### GWDG gitlab

- Web based versioning system
- Advice and support for the establishment of your projects by GWDG
- Connection to the GWDG user administration
- Central monitoring, system stability and backup by the GWDG
- Carpentries workshops on gitlab by SUB





# Persistent identifiers (PIDs)



404 Not Found – 👧

Die angeforderte Seite/Datei konnte nicht gefunden werden.

#### KONTAKT

Georg-August-Universität Göttingen Wilhelmsplatz 1 37073 Göttingen Tel. 449 551 39-0

#### ONLINE-DIEN

Vorfesungsverzeichnis und Personensuche (UniV2) Prüfungsverwaitung (FlexNw Lernmanagement (Stud.IP) Studierendenportal (eCampus Mitarbeiterinnenund Mitarbeiterportal (MaP) Stellenwurst Göttnigen Stellenwurst Göttnigen

#### SERVICE Datenschutz Kontakt Notfall Lageplan

16.06.2020

Alliance

#### eResearch What are persistent identifiers (PIDs)?

Common references like URLs point directly to the location of an object:

If the location of the object changes, the reference points to nothing. The referenced object cannot be found:

The basic idea behind the concept of PIDs is to introduce an **intermediary** between the reference and the referenced object:



This agent monitors all movements or changes made to the object and always forwards requests to the current location.

Adapted from: Kálmán, Tibor: <u>Sustainable referencing of digital objects using persistent identifiers</u>, nestor/DigCurV School 2012, 22-24 October 2012. 16.06.2020 5:

lliance

#### eResearch What are persistent identifiers (PIDs)?

- Prevention of dead links
- Unique naming (referencing) of a digital resource (e.g. journal article or research data)
- Assignment of a permanent and uniquely referenceable code to be resolved on the internet
- Examples:
  - DOI 10.17192/bfdm.20181.7816
  - Handle hdl:11304/6eacaa76-c275-11e4-ac7e-860aa0063d1f
  - EPIC 21.11101/0000-0000-9D43-4
  - URN urn:isbn:0451450523
  - PURL http://purl.abcd.org/ABC/DEF/200
  - ORCID 0000-0001-2345-6789

Alliance



## How do I get a PID?

- Publications: From the publishing journal or repository
- Data: Depositing in a data repository
- At the Göttingen campus:
  - GRO.data (Campus repository)
  - SUB Göttingen: GOEDOC, GoeScholar, University Press
  - GWDG: ePIC PID Service



# **Describing data**



#### Explain these data

										U /
1000025	5	1	1	1	2	1	3	1	1	2
1002945	5	4	4	5	7	10	3	2	1	2
1015425	3	1	1	1	2	2	3	1	1	2
1016277	6	8	8	1	3	4	3	7	1	2
1017023	4	1	1	3	2	1	3	1	1	2
1017122	8	10	10	8	7	10	9	7	1	4
1018099	1	1	1	1	2	10	3	1	1	2
1018561	2	1	2	1	2	1	3	1	1	2
1033078	2	1	1	1	2	1	1	1	5	2
1033078	4	2	1	1	2	1	2	1	1	2
1035283	1	1	1	1	1	1	3	1	1	2
1036172	2	1	1	1	2	1	2	1	1	2
1041801	5	3	3	3	2	3	4	4	1	4
1043999	1	1	1	1	2	3	3	1	1	2
1044572	8	7	5	10	7	9	5	5	4	4
1047630	7	4	6	4	6	1	4	3	1	4



#### Explain these data

Sample code number	Clump Thickness	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Bland Chromatin	Normal Nucleoli	Mitoses	Class
										(2 for benign, 4 for
id number	1-10	1-10	1-10	1-10	1-10	1-10	1-10	1-10	1-10	malignant)
1000025	5	1	1	1	2	1	3	1	1	2
1002945	5	4	4	5	7	10	3	2	1	2
1015425	3	1	1	1	2	2	3	1	1	2
1016277	6	8	8	1	3	4	3	7	1	2
1017023	4	1	1	3	2	1	3	1	1	2
1017122	8	10	10	8	7	10	9	7	1	4
1018099	1	1	1	1	2	10	3	1	1	2
1018561	2	1	2	1	2	1	3	1	1	2
1033078	2	1	1	1	2	1	1	1	5	2
1033078	4	2	1	1	2	1	2	1	1	2
1035283	1	1	1	1	1	1	3	1	1	2
1036172	2	1	1	1	2	1	2	1	1	2
1041801	5	3	3	3	2	3	4	4	1	4
1043999	1	1	1	1	2	3	3	1	1	2
1044572	8	7	5	10	7	9	5	5	4	4
1047630	7	4	6	4	6	1	4	3	1	4

16.06.2020

**Source:** <u>Breast Cancer Wisconsin (Diagnostic) Data Set</u>, Dua, D. and Graff, C. (2019). <u>UCI Machine</u> <u>Learning Repository</u>. Irvine, CA: University of California, School of Information and Computer Science.

### What is metadata?

Metadata are data "descriptions"

- WHO recorded the data?
- WHAT is the content of the data?
- WHEN were the data recorded?
- WHERE (geographically) were the data collected?
- HOW were the data recorded?
- WHY were the data recorded?

58



Source: USFWS - Pacific Region on flickr





### What is metadata?



- Different definitions, depending on the perspective
- Practical approach: Metadata...
  - describe objects in a *structured and standardized* way
  - can help in the selection and identification of resources
  - can describe how to use data correctly or how to reproduce it
  - can describe everything: literature, a painting, places, a set of data...
  - can be digitally linked to objects (embedded) or added separately



### Metadata is everywhere

п

01



Show record						unt per Servi	ng			
					ries 90	Calories fro	m Fat 20			
				2 Files			% Dai	ly Value*		
				File name	RE-DATA-Season1 2017-03-02.xlsx	l Fat 2g		3%		
				Creation date	2017-04-24 21:00:41	aturated Fat 1	5g		8%	
				Lest and Bestler date	2017-04-24-21:03:41	esterol 90mg			30%	
				Last modification date	um 85ma	4%				
				Content type	application/vnd.openxmlformats-officedocument.spreadsheetml.sheet	I Carbohydrate 2g			1%	
		File size	625848	ietary Fiber Og			0%			
				Checksum	48928b915f2e06eafa7d9105c3a7132b	ugars 1g				
				Checksum algorithm	MD5	<b>ein</b> 15g			30%	
	and a start of the	EXIF		CDSTAR ID	EAEA0-0FEC-9E6F-7147-0/RE-DATA-Season1_2017-03-02.xlsx	nin A	0% •	Vitamin C	0%	
F Editor		Download Delete file	ent daily Values are based on a 2,000 calorie diet. Your daily es may be higher or lower depending on your calorie needs:							
10			27 04 2017				Calories	2,000	2,500	
			2110 112011	File name	RE-season1-evaluation_2017-03-02.docx	at	Less than	65g	80g	
			40	Creation date	2017-04-24 21:09:40	sterol	Less than	300mg	200mg	
100	£11.0	1 /005 -		Last modification date	2017-04-24 21:09:40	m sium arbohydrates arv Fiber	Less than Less than	2,400mg 3,500mg	2,400mg 3,500mg	
100	100 7/1.8	1/285 S	+0.0 EV	Content type	application/vnd.openxml formats-officed ocument.word processing ml.document			300g 25g	375g 30g	
		<u>√</u> 7		File size	2824324	n		50g	65g	
-   [®_]		汉	RGB	Checksum	f9cf1a2aa9ab55865b71aea4f7300e59			Ø	PUBLIC	
				Checksum algorithm	MD5				JOIMAIN	
	Autor			CDSTAR ID	EAEA0-0FEC-9E6F-7147-0/RE-season1-evaluation_2017-03-02.docx					
Cop	yright			Download Delete file						

#### 16.06.2020

▼ EX

ISO

#### Metadata: What needs to be in there?

how,

Who created what,



Timo Gnadt gnadt@sub.unigoettingen.de

Contents:

#### abs 36 535 0.5 2 777 2.6 67 1.3 61

Excel table with test data for training purposes

when,



July 26, 2016



On my Office Windows PC

For use in a workshop

eResearch

where and why?

#### Random number generator for changing field data

#### **Description of** the object •

- Methodology and instruments •
- **Units** of measurement
- **References** to related data
- **Definitions of** jargon, acronyms, code
- **Technical information** about the file

Alliance

## Why Metadata Standards?

- A standard provides a structure with which data can be described:
  - Common terms to ensure consistency
  - Common definitions for easier interpretation
  - Common language to facilitate communication
  - Common structure for quick information retrieval
- For search and retrieval, standards offer:
  - a documentation structure in a reliable and predictable format for computer interpretation
  - a uniform summary description of the data set



lliance

eResearch



#### Metadata standards



#### Seeing Standards: A Visualization of the Metadata Universe

Jenn Riley, Devin Becker http://jennriley.com/metadatamap/





#### Metadata standards



Seeing Standards: A Visualization of the Metadata Universe

Jenn Riley, Devin Becker http://jennriley.com/metadatamap/





# **FAIR Data Principles**

### FAIR Data Principles

#### Compilation of guiding principles for research data

- Objective: To make data findable, accessible, interoperable and reusable<sup>\*</sup>.
- Address data producers and data publishers to promote the maximum use of research data
- FAIR refers to both people and machines
- Published in 2016:

Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci Data 3, 160018 (2016). https://doi.org/10.1038/sdata.2016.18

Available here:

https://www.force11.org/group/fairgroup/fairprinciples







# Sharing data

## Data sharing - Motivation



Citation from: William E. Demming (1900-1993)

Alliance

eResearch



### Why share?



Source: Sharing by ryancr via flickr



## Why share?

#### Reputation

- Receive recognition for high quality research
- Better understanding of your methods
- Enables others to review the work
- Recognition for contributions to the research community
- Expand research beyond your discipline



#### Funding

- The provision of data and/or publications may be a requirement of your funding agency.
- It can make your project application more attractive, even if sharing data is not required.



### eResearch

## Why share?



#### Impact

- Sharing makes your data:
  - Easier to find
  - More easily accessible
- Open data/publications lead to increased citations

Source: Richard Matthews on flickr: dart (2011)

#### Reuse

- Starting point for a supplementary study
- Test data for new software and algorithms
- Data for teaching purposes
- Contexts that cannot be assessed at present





#### Data sharing - concerns



- Stockpiling for bad times
- Nobody likes to polish
- Dirt behind the scenes
- Atmosphere of fear
- Small fishes & unicorns

**Source:** <u>All he does is ea</u> Jannes Pockele via flickr Value over time? Embargo! Self use Do it for No documentation Yourself Work in progress "Working set Theft and abuse Trust in science Insignificance The future is unpredictable


## Data sharing - Real barriers

- Place
  - No common tradition
  - No repository
  - No expertise
- Funds
  - no funds
- Rights
  - No carte blanche



Source: <u>Simatai Great Wall</u> by Arian Zwegers on Wikimedia Commons



## Modes of sharing

## Way of transfer

Access conditions

Usage condition

Peer-to-peer Webspace repository restricted upon request Embargo open

none Agreement

License

# Finding Open Access journals and repositories

Alliance

75



# GRO.data: Research Data Repository

GÖTTINO	GEN <b>R</b> ESEARCH <b>O</b>	NLINE			Search -	User Guide	Support	Log in
II Metrics	622 Downloads						Contact (	3 Share
Publish y	our research data	! Search, find,	and cite data from the	e Göttingen Cam	pus an	d beyond.		
Göttingen Rese group of SUB a	earch Online is an institution and GWDG. If you are interes	al repository for the pu sted in publishing your	blication of research data at the Gö data here, please see our author ins	ttingen Campus. It is managet in touch	ged by the with us. See	Göttingen eRese e our Quick Star	arch Alliance t Guide	, a joint
Search this dataverse		Q Find Advanced Search					+ Add Data	
		1 to 10 of 33 Result	s				11	Sort -
		Image review on mobile devices for suspected stroke patients: Evaluation of the mRay® software solution May 14, 2019 - Alex Brehm Dataverse						
		Brehm, A https://do Supplementary data	software solu	tion", ion				
		Alex Brehm Dataverse (Dienste der GWDG)						8
Author Name Hoffmann, Ellen (5)		May 3, 2019 - GRACE Roertgen, Steffen; Harald, Kusch; Claudia Engelhardt; Sven Bingert; Valeria Savin; Inga Kraus, 2019, "PDF Copy Of Online Survey", https://doi.org/10.25625/R48ZD5, Göttingen Research Online, V1						
Elsner, Ines (3) Inga Kraus (3)		This is a PDF copy of an online-survey sent to members of UMG within the GRAcE-project. Link to a copy of the online-sur Survey Survey results Mar 20, 2019 - Dehradun Dataverse Hoffmann, Ellen, 2019, "Survey results", https://doi.org/10.25625/OTNSMI, Göttingen Research Online, V1, UNF-6:dgChlu7R9+FFug99x8e9Q== [fileUNF]						e
Kusch, Harald (3)	) More							
Subject Medicine, Health	th and Life Sciences (14) nities (7) ences (6)							
Arts and Humani Agricultural Scier		Results of the survey of 100 migrant households in Dehradun, Uttarakhand, Northern India (data digitalized into Excel)						
Physics (2)		Kandai-images Mar 20, 2019 - Dehradun Dataverse						
	More	Hoffmann						
Keyword Term Applied Art (3)		Time serie	es of land use maps of Kandai village, U	Ittarakhand, North India				

Alliance

## zenodo

Search

Upload Communities

Q

### Recent uploads

September 1, 2017 (v20) Software Open Access

#### matplotlib/matplotlib v2.1.0rc1

Michael Droettboom; Thomas A Caswell; John Hunter; Eric Firing; Jens Hedegaard Nielsen; Nelle Varoguaux; Benjamin Root; Elliott Sales de Andrade; Phil Elson; Darren Dale; Jae-Joon Lee; Jouni K. Seppänen; Antony Lee; Ryan May; Damon McDougall; David Stansby; Andrew Straw; Paul Hobson; Tony S Yu; Eric Ma; Christoph Gohlke; Steven Silvester; Charlie Moad; Adrien F. Vincent; Jan Schulz; Peter Würtz; Federico Ariza; Cimarron; Thomas Hisch; Nikita Kniazev

#### matplotlib: plotting with Python

Uploaded on September 1, 2017 19 more version(s) exist for this record

#### August 30, 2017 (v1) Working paper Open Access

#### Introducing Parsl: A Python Parallel Scripting Library

Babuii, Yadu; Brizius, Alison; Chard, Kyle; Foster, Ian; Katz, Daniel S.; Wilde, Michael; Wozniak, Justin

Researchers frequently rely on large-scale and domain-specific workflows to conduct their science. These workflows may integrate a variety of independent software functions and external applications. However, developing and executing such workflows can be difficult, requiring complex...

Uploaded on August 30, 2017

#### August 24, 2017 (v2) Dataset Open Access

View

View

#### Aligned ISNI and Ringgold identifiers for institutions

Delpeuch, Antonin

This dataset provides a correspondence between ISNI and Ringgold identifiers, by combining two datasets: Open ISNI for Institutions, available at http://isni.ringgold.com/, which provides metadata for institutions identified by ISNI. The dataset of institutions used by ORCID for disambiguation,...

Uploaded on August 24, 2017 1 more version(s) exist for this record

August 22, 2017 (v2) Dataset Open Access

View

Supplementary Data: Status of the scalar singlet dark matter model (arXiv:1705.07931)

View

#### Zenodo now supports DOI versionina!



Read more about it, in our newest blog post.

#### Using GitHub?

repositories.

Just Log in with your GitHub account and

#### Zenodo in a nutshell

click here to start preserving your

- Research. Shared. all research outputs from across all fields of research are welcome! Sciences and Humanities, really!
- Citeable. Discoverable. uploads gets a Digital Object Identifier (DOI) to make them easily and uniquely citeable.
- Communities create and curate your own community for a workshop, project, department, journal, into which you can accept or reject uploads. Your own complete digital repository!
- Funding identify grants, integrated in reporting lines for research funded by the European Commission via OpenAIRE.
- Flexible licensing because not everything is under Creative Commons.
- Safe your research output is stored safely for the future in the same cloud infrastructure as CERN's own LHC research data

Read more about Zenodo and its features.

## Other services on campus

Name	available through	Purpose/comments
Jupyter notebooks	GWDG	Live editing and execution of text, diagrams, equations and code in a web browser
CodiMD pad	GWDG	Collaborative text editing
Electronic laboratory notebook	UMG	(Re)organisable, searchable and storable research documentation
Self-study online courses on digital competencies	SUB	Courses on literature search, IT basics, data security, data visualization, OER <i>https://www.uni-goettingen.de/en/565228.html</i>
Open Access Publication Fund	SUB	full coverage for up to € 2.000,- for publication in OA journals
Video Conferencing	GWDG	Various solutions

Alliance

eResearch



# **GWDG** services

## SERVICES

#### Storage Services

Data Archiving Backup GWDG CrashPlan PROe File Service GWDG Cloud Share GWDG ownCloud Cryptshare

#### Application Services

Persistent Identifier (PID) High Performance Computing Library Service Aleph Database Service Oracle Database Service MySQL Application and Registration Services Plagiarism Prevention Online Surveys Bioinformatics Programs Statistics Programs Jupyter

#### Email & Collaboration

E-Mail-Service (MS Exchange 2010) Spam and Virus Filtering Mailing Lists MS Sharepoint Managed Services Project Management Service Etherpad ShareLaTeX Rocket Chat GitLab

#### Server Services

Virtual Server Hosting/Housing of Servers Web Hosting GWDG Cloud Server FTP-Server Puppetserver

#### Network Servies

IP Address Management System Cable and Route Management System System Monitoring Setting up eduroam Integration into the Active Directory User Management with OpenLDAP Client Management Client Management for macOS and iOS

#### IT Security Services

Vulnerability Scans on Network-attached Equipment Public-Key-Infrastructure (PKI) Authentication and Authorization Infrastructure (AAI) Virus Protection (Sophos Update Service)

#### General Services

Identity Management Courses Software and Licence Management Videoconferencing Computer Lending Pool Print & Scan Services

#### IT Consulting

Scientific Data Management IT Security Hardware Purchase Apple Support Centre Establishing Directory Services (AD, LDAP) Planning of Data Transmission Networks

### https://www.gwdg.de/services

# Further reading

eResearch

- DataOne
  - <u>https://www.dataone.org/education-modules</u>
- MANTRA
  - <u>http://mantra.edina.ac.uk/</u>
- "Reproducible research with jupyter notebooks"
  - <u>https://reproducible-science-curriculum.github.io/rr-jupyter-workshop/</u>
- "IT in a nutshell" a self-learning course by SUB Göttingen on IT basics
  - https://www.uni-goettingen.de/en/613871.html

# Göttingen eResearch Alliance

## Team

- Various discipline-specific expertise
  - mainly in the natural sciences, humanities, computer science
- Jointly operated by



SUB | NIEDERSÄCHSISCHE STAATS- UND UNIVERSITÄTSBIBLIOTHEK GÖTTINGEN

- Partner: Department of Research, University Medicine (UMG)
- Comprehensive expertise on eResearch topics



eResearch

Alliance

# What the eRA can do for you

- Consulting / Support / Services
  - Research Data Management
  - Publication strategies
  - Digital methods, software and technologies to improve research projects
  - Central point of contact for experts & expertise throughout the campus

- Training
  - Discipline- or projectspecific or general

eResearc

- Information material / knowledge database
- Cooperation
  - Project partnership
  - Project as a service

## www.eresearch.uni-goettingen.de

lliance



# Thank you for your participation!

The comments, questions and suggestions of the participants for this training can be found here:

https://pad.gwdg.de/LoveYourData\_20200812?both

## CONTACT:

info@eresearch.uni-goettingen.de

www.eresearch.uni-goettingen.de



## Traceable:

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

F2.data are described with extensive metadata.

F3 (Meta) data are registered or indexed in a searchable resource.

F4.metadata specify the data identifier.



## Accessible:

A1. (Meta-)data can be <u>retrieved by</u> <u>means of their identifier</u> using a standardised <u>communication protocol</u>.A1.1 the protocol is open, free and <u>universally implementable</u> .A1.2 the protocol allows an <u>authentication and</u> authorisation procedure if required.

A2.metadata are <u>accessible</u> even if the data are no longer available.



## Interoperable:

I1. (Meta-)data use a formal, accessible, shared and generally applicable <u>language to represent</u> knowledge

I2. (Meta-)data use <u>vocabularies</u> that follow the FAIR principles.

I3 (Meta) data contain <u>qualified references</u> to other (meta) data.



## Reusable:

R1.meta(data) have a variety of accurate and relevant <u>attributes</u>

.R1.1. (Meta)data are released with a clear and accessible <u>data use licence</u>.R1.2. (Meta)data are associated with their <u>provenance</u>.R1.3. (Meta)data meet domainrelevant <u>Community standards</u>.