# Metadata and controlled vocabularies – a treasure map to your data

**Oliver Biehlmaier**
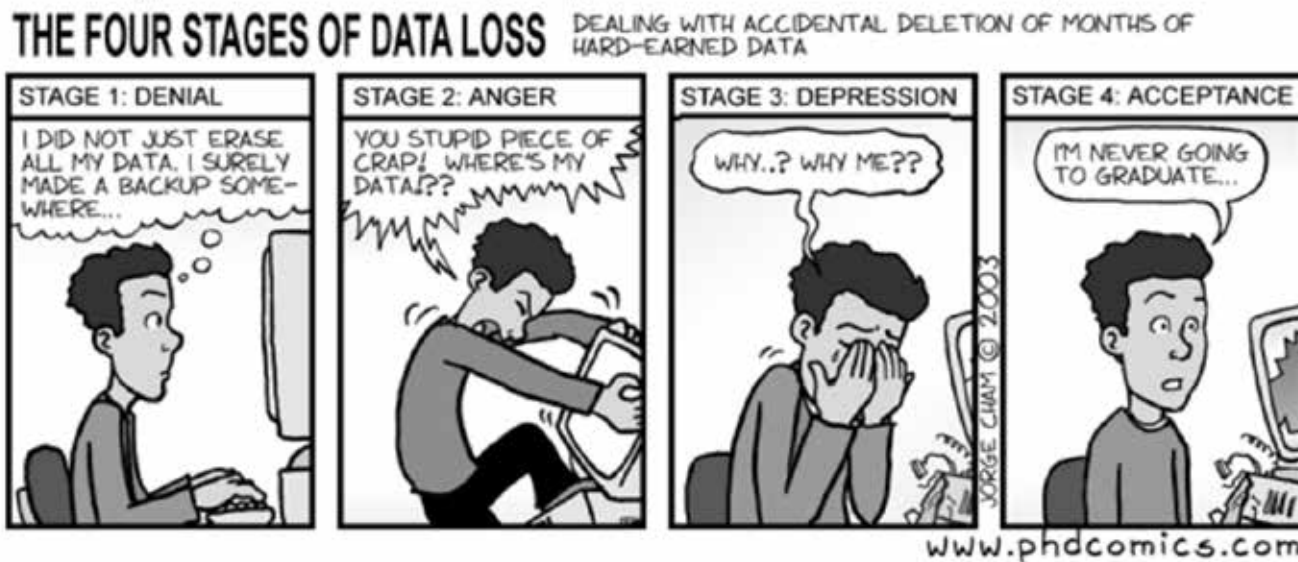
Imaging Core Facility, Biozentrum, University of Basel

# Why to manage and share your data

- You act **according to funding agencies requests**



**Piled Higher and Deeper** *by Jorge Cham*          www.phdcomics.com

THE FOUR STAGES OF DATA LOSS   DEALING WITH ACCIDENTAL DELETION OF MONTHS OF HARD-EARNED DATA

STAGE 1: DENIAL — I DID NOT JUST ERASE ALL MY DATA. I SURELY MADE A BACKUP SOMEWHERE...

STAGE 2: ANGER — YOU STUPID PIECE OF CRAP! WHERE'S MY DATA!??

STAGE 3: DEPRESSION — WHY..? WHY ME??

STAGE 4: ACCEPTANCE — I'M NEVER GOING TO GRADUATE...

JORGE CHAM © 2003

www.phdcomics.com

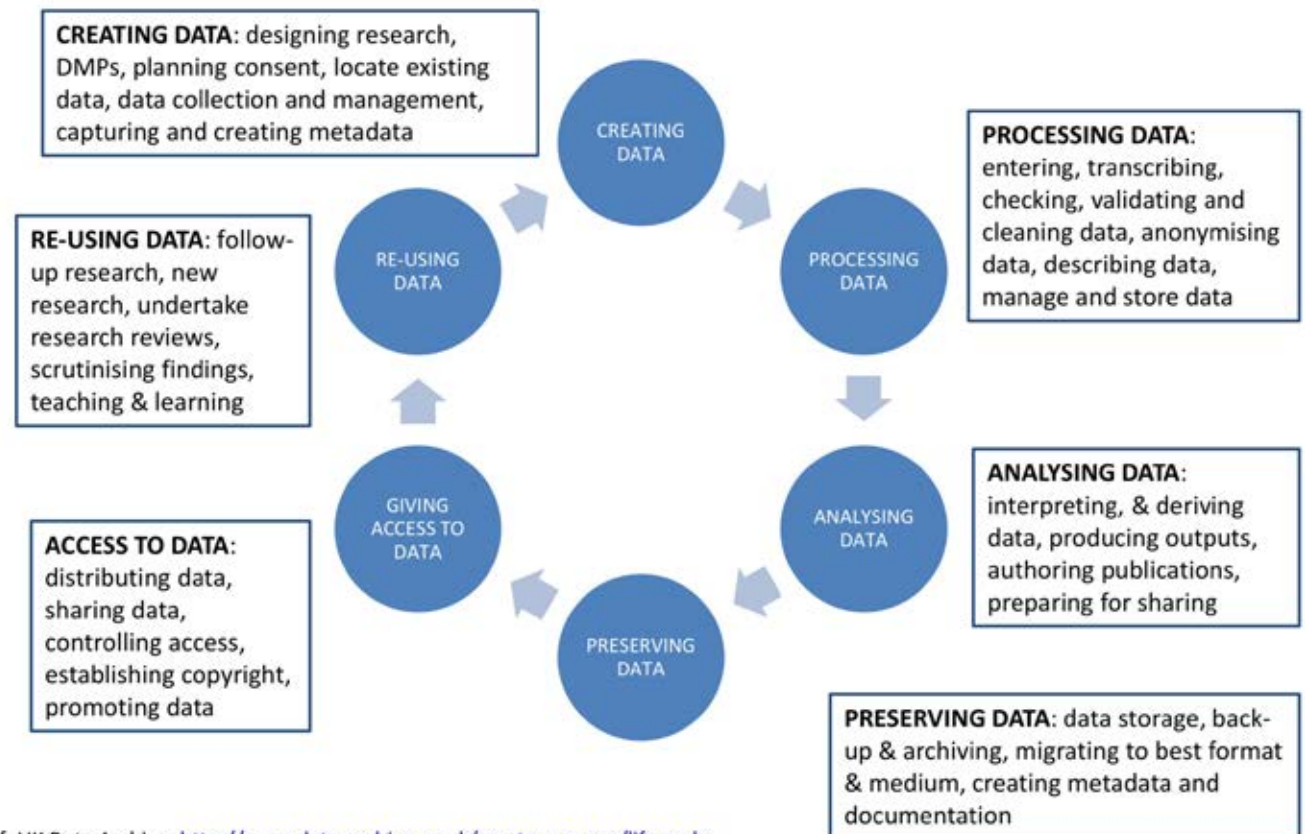title: "Stages of Data Loss" - originally published 10/30/2003

# Why to manage and share your data

- You act **according to funding agencies requests (SNF, ERC)**
  **=>that's where the actual momentum comes from!**

- **Find and understand your data when you need to use it**

- **Continuity if project staff leave or new researchers join**

- **Avoid unnecessary duplication** *(e.g. re-collecting or re-working data on different storage shares)*

- **Data underlying publications are maintained**, allowing for validation of results

- Enables **more collaboration and advances research**

# How/where to start? Where to go?

- Think about your data creation process and beyond!



## Research data lifecycle

**CREATING DATA**: designing research, DMPs, planning consent, locate existing data, data collection and management, capturing and creating metadata

**PROCESSING DATA**: entering, transcribing, checking, validating and cleaning data, anonymising data, describing data, manage and store data

**ANALYSING DATA**: interpreting, & deriving data, producing outputs, authoring publications, preparing for sharing

**PRESERVING DATA**: data storage, back-up & archiving, migrating to best format & medium, creating metadata and documentation

**ACCESS TO DATA**: distributing data, sharing data, controlling access, establishing copyright, promoting data

**RE-USING DATA**: follow-up research, new research, undertake research reviews, scrutinising findings, teaching & learning

Ref: UK Data Archive: http://www.data-archive.ac.uk/create-manage/life-cycle

# Lifecycle in bioimaging



- Choice of imaging & staining methods
- Instrument choice and setup
- Novel imaging technique development?
- Plan reproducible analysis workflows
- Consider availability of re-usable datasets?

Planning, Data Mining & Re-use

Image Acquisition
- Quality control of imaging setups and acquisition settings
- Which metadata is captured automatically (technical metadata)?
- Documentation of sample preparation workflows & protocols
- Which file formats are generated?
- Enrich metadata with editation tools?

- Trusted public or institutional repository?
- Upload to added-value database?
- Choose data and metadata license
- Choose suitable data format
- Metadata enrichment before upload

Archiving, Long-Term Storage

FAIR

Image File Storage & Access
- What local storage infrastructure is available?
- Is data stored temporarily, does it have to be transfered?
- Is an imaging data management software available?
- Necessity of large file transfer & storage?
- How is data secturity (avoiding data loss, data privacy) managed?
- Is cloud computing & -storage required?

- Check open access requirements
- Check publisher policies
- Check funder policies
- Adhere to figure composing standards
- Consider dataset publication & preprints
- Link publication and data (persistent identifiers)

(Data-) Publication

Processing & Bioimage Analysis
- Are shared digital environments required?
- Choice of analysis tools and software
- Use of (semi-)automated analysis workflows
- Use of AI-based analysis tools?
- Ensure analysis data provenance
- Linking original and derived data

*from I3D:bio – Information Infrastructure for BioImage Data*

# Two most important points to care about when planning data management

1. Make sure you provide all possible **Metadata**

NO METADATA NO FUTURE

2. Be **F.A.I.R.**

Findable Accessible Interoperable Reusable

# What is METADATA?

## "We kill people based on metadata"

Gen. Michael Hayden
Former head of the National Security Agency

# What is METADATA?

**Data about your data...**



YO DAWG, I HEARD YOU LIKE METADATA

SO I MADE SOME METADATA ABOUT YOUR METADATA ABOUT YOUR DATA

# What is METADATA?

**METADATA = a set of data that describes and gives information about other data.**

- **scientific images** – image size, objective (mag., NA), filters, laser wavelength, exposure time, opt. slice thickness, stack size, xy resolution, etc.

- **photography** – image size, exposure, time & date, objective, zoom, flash/no flash, GPS data, etc. (EXIF information)

- **movies** – file format, movie length, size, actors, director, producer, etc.

- **ebooks** – format, author, editor, year of publication, # pages, ISBN, etc.

- **Music streaming** – artist, album, title, length of song, genre, etc.

# What is METADATA? – a definition

"Metadata is constructed, constructive, and actionable."

*Karen Coyle, Digital Librarian, Author of Coyle's InFormation*

- Constructed - a man-made artifice, not naturally occurring
- Constructive - serving a useful purpose, to solve some problem
- Actionable - can be acted upon, processed by humans and machines

# METADATA types

**There are 3 main types of metadata:**

- **Descriptive metadata** enables **discovery**, **identification**, and **selection** of resources. It can include elements such as title, author, and subjects.

- **Administrative metadata** facilities the **management of resources**. It can include elements such as technical, preservation, rights, and use.

- **Structural metadata,** generally **used in machine processing**, describes **relationships** among various parts of a resource, such as chapters in a book.

# METADATA in bioimaging (example)

- **Technical metadata** contains information on hardware and settings used to acquire an image.
  - includes **device specifications**, **objective** lens specifications, **light-source**, **laser** and **filter** settings, number of **channels**, **camera**, bit depth, etc.
  - is **automatically recorded** by most microscopes and **stored in the metadata header** of the image file.
  - can be accessed and edited using metadata editing tools.

- **Experimental and sample preparation metadata** contains information about the specimen that is imaged (organism, cell line, organ, expression contstructs, etc.)
  - The **experimenter must add** all this information to the data.
  - This experimental metadata **should also cover**
    - **how a sample has been prepared** for imaging
      - fixed and stained with antibodies
      - conjugated to fluorescent dyes?
    - whether **live** cells and under what conditions were they imaged?

# METADATA you need in bioimaging - visualization



from REMBI: Recommended Metadata for Biological Images—enabling reuse of microscopy data in biology

# METADATA in bioimaging



© Thao Do (Allen Institute, Seattle, WA, USA)

Nature Methods FOCUS issue on Reporting and Reproducibility in Microscopy
https://www.nature.com/collections/djiciihhjh

All information that is needed to interpret, evaluate the quality, reproduce and share microscopy images

- Sample preparation
- Image Acquisition
  - Hardware configuration
  - Acquisition setting
  - Quality Control
- Image data processing and analysis

# How to create useful experimental METADATA

You want to describe your experimental data in a F.A.I.R. way?

**What would you do?**

**You should use**

- specific and defined **keywords or tags**
- **Controlled vocabularies** or **taxonomies**

# Treasure map = METADATA and controlled vocabularies

How do you **find your treasure**?
- Randomly checking every island you find?
- Searching every corner of the island and do some digging at suspicious locations?

**No!!!!**

You would only start searching **if you had the treasure map**!
- Directly navigate to the island that fits the map
- Follow the path indicated on the map which is referring to **specific known or easily identifiable landmarks**
- **Search at the location** where the data treasure is hidden

**Metadata is a map. Metadata is a means by which the complexity of an object is represented in a simpler form.**

# Standardizing keywords – an example

- **Star Wars: Episode I -- The Phantom Menace**
- Episode 1
- Episode I
- Phantom Menace
- Star Wars Episode I The Phantom Menace
- Star Wars Episode I: The Phantom Menace
- Star Wars prequel
- Star Wars: Episode 1 -- The Phantom Menace
- Star Wars: Episode i -- the Phantom Menace
- Star Wars: Episode I: The Phantom Menace
- Star Wars: Episode I--The Phantom Menace
- Star Wars: Episode I--The Phantom Menance
- Star Wars: Episode One -- The Phantom Menace
- Star Wars: The Phantom Menace
- Star Wars: The Phantom Menace -- Episode I
- The Phantom Menace
- The Phanton Menace

*from https://www.slideshare.net/rlovinger/metadata-is-a-love-note-to-the-future*

# Controlled vocabularies example: Biological Taxonomy

**Carolus Linnaeus,** first person to combine binomial definition (nomenclature) with a hierarchical structure of classification.

His **system organized both plants and animals from the level of the kingdom, right down to species**.

He used this system consistently to identify every species of plant and animal he came across and this is the basis of the system we use today.



Let's take these 5 animals:

© Shutterstock, 2019.

# Controlled vocabularies example: Biological Taxonomy



©Shutterstock, 2019

# Controlled vocabularies example: Biological Taxonomy

# Controlled vocabularies example in life sciences today - PRIDE

# Medical Subject Headings - MESH

# What if you do not know "your" specific ontology yet?

# Search OLS (Ontology Lookup Service)

# Ontologies that might be useful for NCCR AntiResist

# What if you do not find a suitable ontology?

**Create your own ontology**

1. Determine the **domain and scope** of the ontology.

2. **Consider reusing** existing ontologies. At least to start with

3. Enumerate **important terms**.

4. Define the **classes & class hierarchy**.

5. Define the **properties of classes**.

*from http://www.ksl.stanford.edu/people/dlm/papers/ontology101/* **University of Basel** 25

# If you are planning to create your own ontology

- Try to **avoid individual group level solutions** at all costs, as they are almost never sustainable and F.A.I.R.

- Try to **reach out to the largest possible community** and get them all involved (if possible).

  This usually makes the process longer and more tedious, but in the end it **results in an ontology that is used by everyone** in the field (see PRIDE, PDB, etc).

  "We just do the right thing to the best of our ability and that happens to be F.A.I.R."

# On of the success stories for a database, metadata, ontology

*from Berman et al 2012*

# Workshop - LEGO group exercise
## (5 groups à 5-6 persons)

- **Group into 5 groups à 5-6 persons**

# LEGO group exercise (5 groups à 5-6 persons)

- **Build the structure assigned to the respective building station**
  - **Document your structure** using available tools
  - **Take a picture** of your build (or tutor will take a picture)

- **Disassemble** build

# LEGO group exercise (5 groups à 5-6 persons)

- **Build #1** - Build the structure assigned to the respective building station
  - **Document your structure** using available tools at the station

  - **Take a picture** of your build (or tutor will take a picture)

- **Disassemble** build

- **Build #2** – Rebuild the previous structure using the documentation of the previous group

- **Compare**
- **Discuss your experience**

# LEGO group exercise- discussion of results

➢Did you find this a simple way to document your process?

➢Was there anything you found difficult to capture?

➢Did those replicating the builds find it straightforward to follow?

➢Did you encounter any ambiguity in the instructions?

# We hope you liked the workshop

And always remember: