# Research Data Management

## HPC Day for Physicists
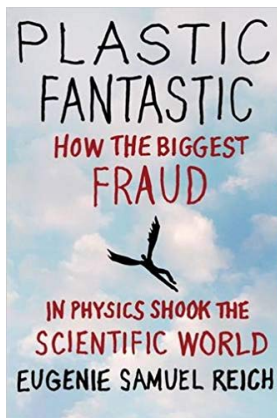


Picture: Harald Kusch, Uni-Medizin Göttingen
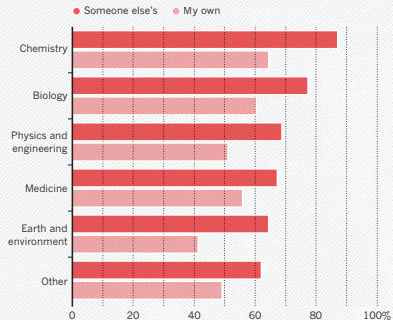
Jörg Steinkamp
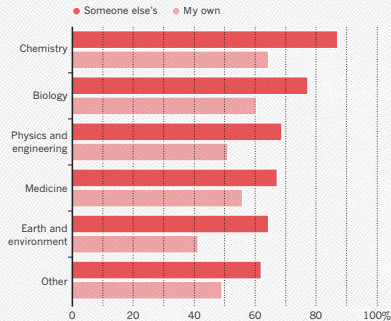
ZDV Zentrum für **Datenverarbeitung**

# Content

# Why care about RDM



HAVE YOU FAILED TO REPRODUCE
AN EXPERIMENT?
Most scientists have experienced failure to reproduce results.

Baker (2016), Nature.

# Why care about RDM



Baker (2016), Nature.



Jorge Cham www.phdcomics.com

# Why care about RDM

Everybody can handle order, but only a genius can
master chaos, ...

# Why care about RDM

Everybody can handle order, but only a genius can master chaos, ...

... but what happens when the genius leaves?!

# Why care about RDM

Everybody can handle order, but only a genius can master chaos, …

… but what happens when the genius leaves?!

# Why care about RDM

## Researchers

- minimize risk of data loss

- maximize ...
    - ... efficiency
    - ... sustainability
    - ... reproducibility

- facilitate ...
    - ... teamwork
    - ... follow-up projects

- increase reputation

- beeing nice to following scientists

# Why care about RDM

## Researchers

- minimize risk of data loss
- maximize ...
    - ... efficiency
    - ... sustainability
    - ... reproducibility
- facilitate ...
    - ... teamwork
    - ... follow-up projects
- increase reputation
- beeing nice to following scientists

## System Administrators

- minimize ...
    - ... risk of data loss
    - ... costs
- maximize ...
    - ... efficiency
    - ... sustainability
- avoid patchwork rug
- provide central on-site storage
- simplify scientific workflows
- beeing nice to all scientists

# Archiving



2. Archiving
   - Backup vs. Archive
   - Money, money, money

3. Data life cycle

4. FAIR principles

5. Tools for RDM

# Backup isn't Archiving and Archiving isn't Backup



|            | **Backup**  | **Archive** |
|------------|-------------|-------------|
| **target** | active data | final data  |
| **purpose** |            |             |
| **timeframe** |          |             |
| **metadata** |           |             |

# Backup isn't Archiving and Archiving isn't Backup

|  | **Backup** | **Archive** |
|---|---|---|
| **target** | active data | final data |
| **purpose** | protection and recovery | preservation of information |
| **timeframe** | | |
| **metadata** | | |

# Backup isn't Archiving and Archiving isn't Backup



| | Backup | Archive |
|---|---|---|
| **target** | active data | final data |
| **purpose** | protection and recovery | preservation of information |
| **timeframe** | short term | long term |
| **metadata** | | |

# Backup isn't Archiving and Archiving isn't Backup



|  | **Backup** | **Archive** |
|---|---|---|
| **target** | active data | final data |
| **purpose** | protection and recovery | preservation of information |
| **timeframe** | short term | long term |
| **metadata** | mainly ACL and file attributes | data-specific |

# Archiving medium



| Medium | logevity | HW cost per PB | Total cost per year |
|--------|----------|----------------|---------------------|
|        |          |                |                     |

# Archiving medium



| Medium | logevity | HW cost per PB | Total cost per year |
|--------|----------|----------------|---------------------|
| **SSD** | 100 yr | 100.000€ | 690.000$ |

# Archiving medium



| Medium | logevity | HW cost per PB | Total cost per year |
|--------|----------|----------------|---------------------|
| **HDD** | 3–5 yr | 32.000€ | 555.000$ |
| **SSD** | 100 yr | 100.000€ | 690.000$ |

# Archiving medium



| Medium | logevity | HW cost per PB | Total cost per year |
|--------|----------|----------------|---------------------|
| **Tape** | 30 yr | 6.000€ | 65.000$ |
| **HDD** | 3–5 yr | 32.000€ | 555.000$ |
| **SSD** | 100 yr | 100.000€ | 690.000$ |

# Archiving medium



| Medium | logevity | HW cost per PB | Total cost per year |
|--------|----------|----------------|---------------------|
| **Tape** | 30 yr | 6.000€ | 65.000$ |
| **HDD** | 3–5 yr | 32.000€ | 555.000$ |
| **SSD** | 100 yr | 100.000€ | 690.000$ |
| **DVD** | 10–25 yr | 60.000€ | – |

**Image stored on a CD-R**

# Archiving: Disc in Desk



**Image stored on a CD-R**

**Still readable, but ...**
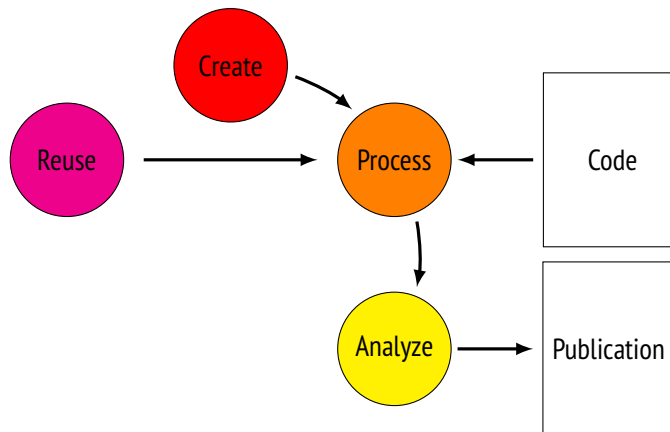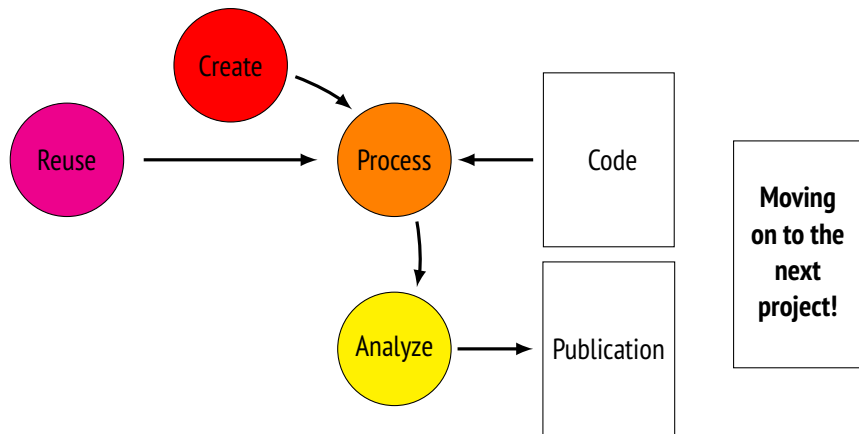
@EdithHalvarsson

# Data life cycle

# Classic data life cycle

# Scientific part of the data life cycle

# Scientific part of the data life cycle
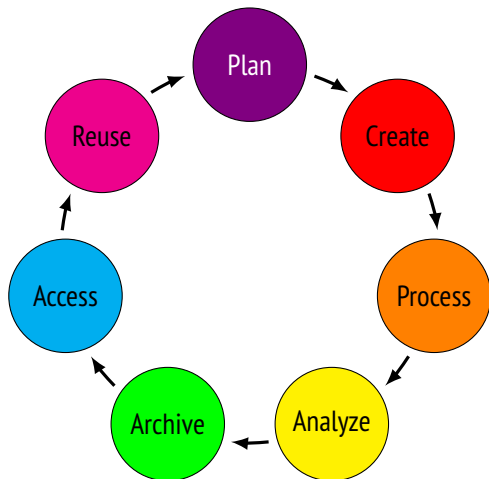
# Scientific part of the data life cycle

# Getting lost in the Research Data Management Cloud



Harald Kusch, Uni-Medizin Göttingen

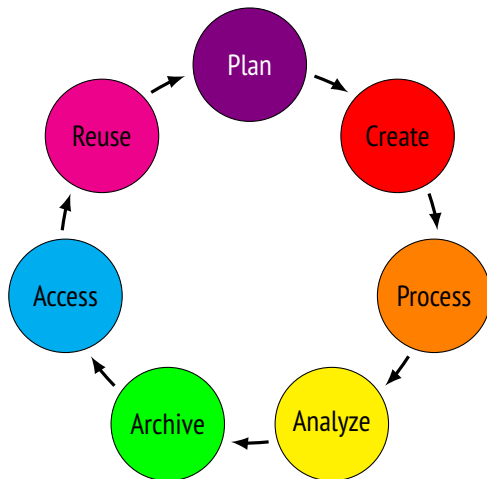# Don't forget the Planning

# Don't forget the Planning

# Don't forget the Planning



## Data Management Plan
- dynamic document

# Don't forget the Planning



## Data Management Plan
- dynamic document
- JGU-RDMO [link]

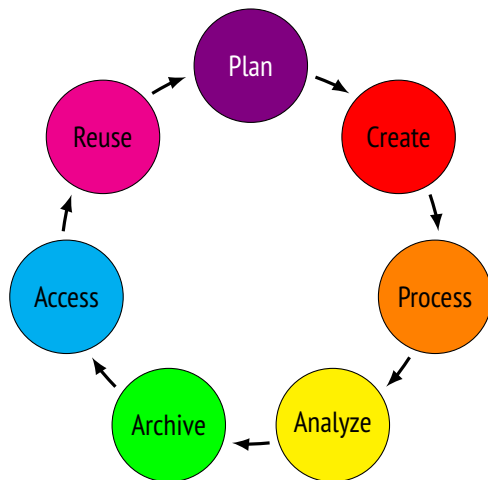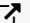# Don't forget the Planning



## Data Management Plan
- dynamic document
- JGU-RDMO 🗗
- other Tools:
  - ▸ DMPonline 🗗
  - ▸ DMPTool 🗗

## Consultation
- HPC: Christian Meesters ✉
- ZDV: Jörg Steinkamp/
  Sarah Wettermann ✉
- Administrativ: Anne Vieten ✉
- Other Universities 🗗

# FAIR principles

gerait@Pixabay.com

Findable  Accessible  Interoperable  Reusable

# **F**indable by keywords

## Sufficient rich metadata

### Minimum set

- Creator
- Title
- Date
- Location
- Publisher
- Keywords
- ...

# **F**indable by keywords

## Sufficient rich metadata

### Minimum set

- Creator
- Title
- Date
- Location
- Publisher
- Keywords
- ...

### Extended metadata

- NFDI

# **F**indable by keywords

## Sufficient rich metadata

### Minimum set

- Creator
- Title
- Date
- Location
- Publisher
- Keywords
- ...

### Extended metadata

- NFDI
- Other/general standards:
  - ▸ RADAR ⌅
  - ▸ Dublin Core ⌅
  - ▸ Data Cite ⌅
  - ▸ Disciplinary Metadata ⌅
  - ▸ Research Data Management Toolkit: Metadata Standards ⌅
  - ▸ Schema.org ⌅
  - ▸ Bioschemas.org ⌅

XKCD

# Let others **A**ccess your (meta-)data



modified after @Cloudways 🔗

©CloudTweaks.com

> ⚠️ **The Cloud is not an archive**
> - ~~Dropbox/Box~~
> - ~~Google Drive~~
> - ~~Owncloud/Nextcloud~~
> - ~~Seafile~~
> - ~~Amazon Drive~~
> - ~~Microsoft Azure~~
> - ...

# 3rd party repositories

## General repositories

-  **DRYAD**
-  fig**share**
-  + zenodo

# 3rd party repositories

## General repositories

-  **DRYAD**
- figshare
- + zenodo

## Lists of subject specific repositories

- Subject specific at
  SCIENTIFIC DATA
- Search engine
  re3data.org
  REGISTRY OF RESEARCH DATA REPOSITORIES
- ...

# Interoperable: Let others be able to read your data

## General

- unencrypted
- (uncompressed)

- open, documented standard
- non-proprietary non-patented

# **I**nteroperable: Let others be able to read your data

## General

- unencrypted
- (uncompressed)

- open, documented standard
- non-proprietary non-patented

## However, ...

- docx, xlsx are acceptable
- raw device data

# **R**eusable: Allow/Enable others to use it

## General

- Provide a license
- add subject-relevant metadata
- describe the source/origin of the data

3 Data life cycle

4 FAIR principles

5 Tools for RDM
- Archiving

# Archiving

# Archiving

## TSM (Tape Library) ⤢

- Two copies at different location
- Encrypted Tapes
- No need to care about ...
    - ... access control
    - ... reuse/-ability
    - ... metadata



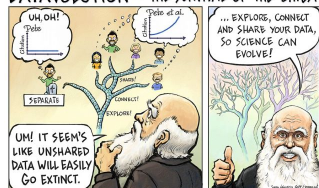publicdomainvectors.org ⤢

# Archiving

## TSM (Tape Library) ↗

- Two copies at different location
- Encrypted Tapes
- No need to care about ...
  - ▶ ... access control
  - ▶ ... reuse/-ability
  - ▶ ... metadata



publicdomainvectors.org ↗

## iRODS ↗

- Two copies at different location
- Encrypted HDDs
- You benefit from ...
  - ▶ ... access control
  - ▶ ... attached metadata
  - ▶ ... publication possibility



Cartoon: Seppo Leinonen, www.seppo.net ↗

# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem

## What is iRODS?

- Virtual Filesystem

# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem

### What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata

# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem ↗

### What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata
- Publishing

# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem $\nearrow$

### What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata
- Publishing
- Fine grained Access Control

# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem ⬈

## What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata
- Publishing
- Fine grained Access Control
- Workflow automation

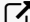# iRODS: **i**ntegrated **R**ule-**O**riented **D**ata **S**ystem ↗

### What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata
- Publishing
- Fine grained Access Control
- Workflow automation
- Commandline tools

### What is iRODS?

- Virtual Filesystem
- Attached searchable Metadata
- Publishing
- Fine grained Access Control
- Workflow automation
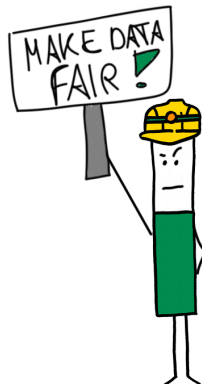- Commandline tools
- (WebUI)

# Upcoming courses and further information

## Further links

- Kompetenzteam FDM 🔗
- Archiving @ Mogon-Docs 🔗

## Courses

- HPC-related courses 🔗
- Git Courses Registration 🔗

# Questions?



Sketchnotes & Sketches – FranziMachtDas ⬀