# Data management intro

Wolfgang Müller, HITS

**HITS**
**Heidelberg Institute for Theoretical Studies**

**Think beyond the limits!**

# FAIRDOM Initiative

**F**indable
**A**ccessible
**I**nteroperable
**R**eusable

**D**ata
**O**perations
**M**odels

http://fair-dom.org

- develop a community
- establish an internationally sustained Data and Model Management service

- Partners from U Manchester, HITS, U Zürich, ETH Zürich
- Funded by BMBF, ISBE, BBSRC, SystemsX (CH)
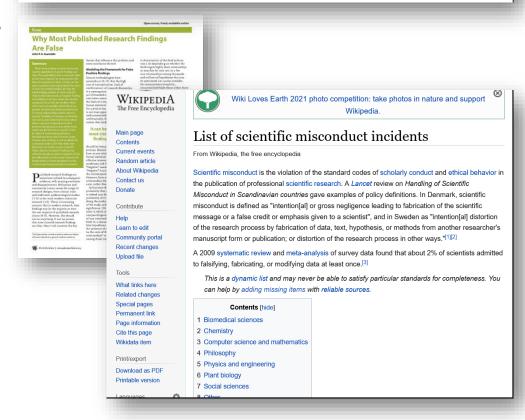- Part of ELIXIR as "Node Resource" in UK and DE

# Topics

- Why data management?

- Some general thoughts leading up to the SEEK day and the openBIS day

HITS · Heidelberg Institute for Theoretical Studies

# The reason: Reproducibility Crisis



- In a variety of subject matters
  - Many positive findings are false

  - Many papers to be trusted
  - But many not to be trusted
  - High profile misconduct incidents

- Need for fast way to access data
  - Proving reproducibility
  - Making control easier
  - Making QC easier within labs

# Consequence: drive for FAIR data

- Next slides:

- What is FAIR

- And what do funders ask for, at the example of DFG

# What is FAIR?
# Data should be like money



To be useful, money must be…

findable



accessible



interoperable



reusable

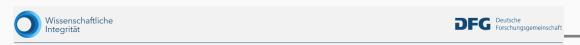# But data is more

| Principle | Meaning |
|---|---|
| **Findable** | Identified, metadata enriched, searchable |
| **Accessible** | Locatable, standardized protocols |
| **Interoperable** | Metadata rich, sustained, metadata survives data |
| **Reusable** | Licensed data&metadata, versatile metadata |

The stick

# Funders want this

# Rough overview of GSP guidelines of DFG

- Comprehensive standards of good research practice

- Many suggestsions with an impact on data management:

  - Quality control
  - Long-term preservation
  - Publishing of research

The carrot

Imagine: You want to discover a project/group/institute

For example:
LiSyM

# Publications

- PubMed: „Liver Lisym" → 13 results

**263** papers
**> 50 Journals**

# Now, let's look at

- Models

- Data

# Models: Model DBs & GitHub/GitLab

# Datacollections, Zenodo, web apps

# Summarising

- Many specialized data sources
- Many useful data locations
- Quite some data merged – curated

- Looking at all of them: too cumbersome

- Data management enables to have
    - **one place to find them all**
    - **one place to standardize**

- Needs to be adapted to
    - Remote data
    - Local data

# Now

- Talk by Piotr Zadora, DKFZ
  Covering the use of openBIS
  in the Klingmüller Department at DKFZ

- Then some more work about challenges
  to address in DM or DMM

- Then „the SEEK day"

- Followed by „the openBIS day"

Some more general words about DM

- Strong push by lead

- Onboarding by proficient users

- Jupyter Notebooks in openBIS

- Data models
(FAIRDOM support but scientific tech organiser, internally, too)

# Goal:
# Project Data Management



Organisation
Communication
Dissemination

Partners
Funders
Public

500 miles

Current, comprehensive, FAIR

# Challenges to adress

- Social & organisational
- Conceptual
  - Data structures
  - Ontologies
  - Standards
- Tooling
  - Which tools to use
  - How to combine

# Assembling data management processes

- Collect raw and processed (secondary) data, models & metadata together with experimental context
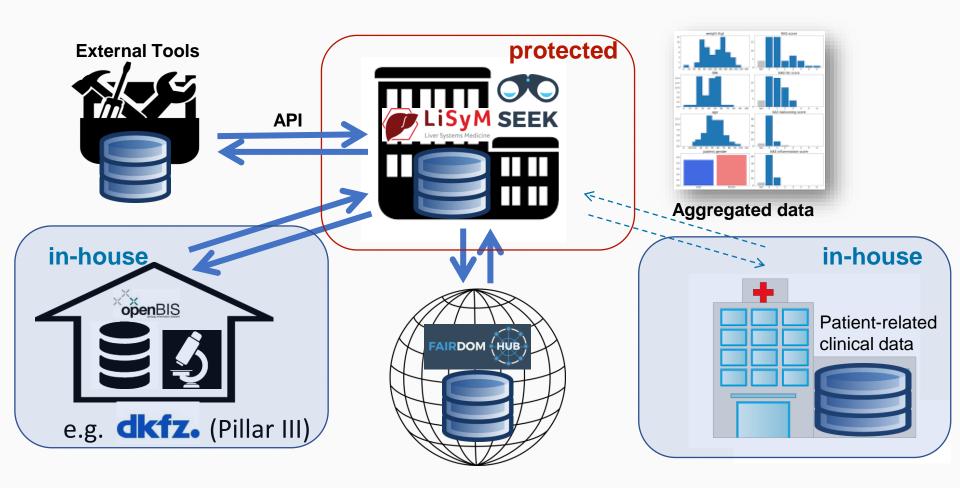
- Organise and link assets

- Prepare reproducible publications

- Use standardised metadata

- Share with colleagues and public

- Integrate with legacy, home grown, external systems

- Reuse tools and community archives



**RDMkit**
https://rdmkit.elixir-europe.org

HITS · Heidelberg Institute for Theoretical Studies

# Structuring deployment: LiSyM Example



External Tools

API

protected

Aggregated data

in-house

openBIS

e.g. dkfz. (Pillar III)

FAIRDOM HUB

in-house

Patient-related clinical data

# Assembling vs. Integrating: integrated JWS modelling tool in SEEK-based FAIRDOMHub



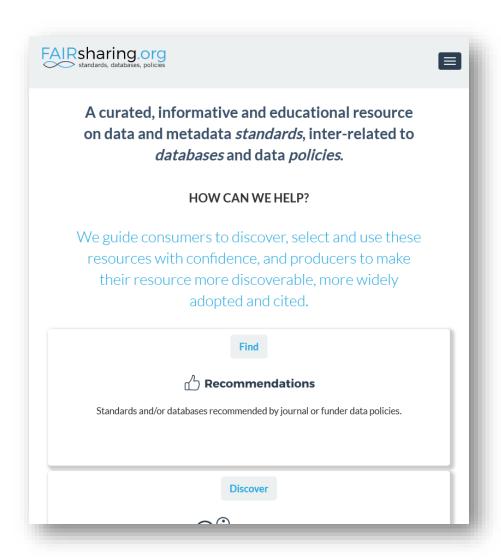https://fairdomhub.org/models/284/simulate?version=1&constraint_based=0

# Ontologies and standards: FAIRsharing.org

**How to structure your data?**
Reuse other people's work

- Community standards

  - 1532 Standards
    (e.g. the MIBBI standards)

  - 1773 Databases

  - 140 Policies

# Social

- **Insight:** Data management means **change**

- **Challenge: How to include people?**
    - How to get information
    - Useful suggestions how to do better
    - Refinement of workflows
    - **Strive for painless ingestion and useful outcome**

- **Proposal:** Have responsibles on per-project basis

- Regular meetings with data management responsibles

- Who should be responsible?
  Depends on **your** organisation

The FAIRDOM approach: Have „PALs", a focus group in the project



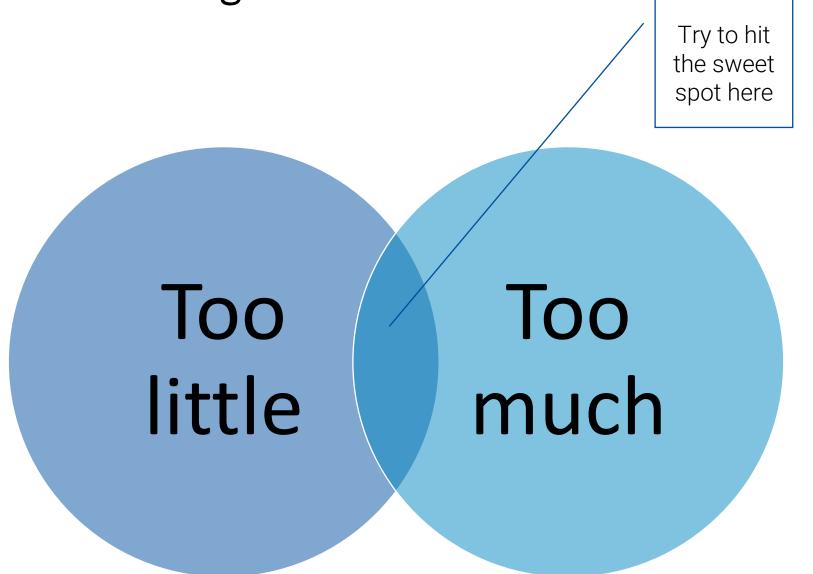PALs kickoff of FAIRDOM 2013

# Organisational: Data sharing policy

**Example: LiSyM DSP, derived from ERASysAPP DSP**

- Expression of intent

- Who is responsible?

- Who shares with whom?

- When?

- What?

- What in case of disagreement?

# Spend the right amount of effort on this

Too little

Too much

Try to hit the sweet spot here

# Now: Front End
## Find, Access and Organise assets

FAIRDOM SEEK

– Upload data
– Link to data regardless of physical store
– Sharing
– ISA structure
– Yellow pages and collaboration
– Supplementary data for publications
– Standards-compliant

# Thanks, Questions?

Wolfgang.mueller@H-ITS.org