# Open Science in all its Facets with a Focus on Research Software
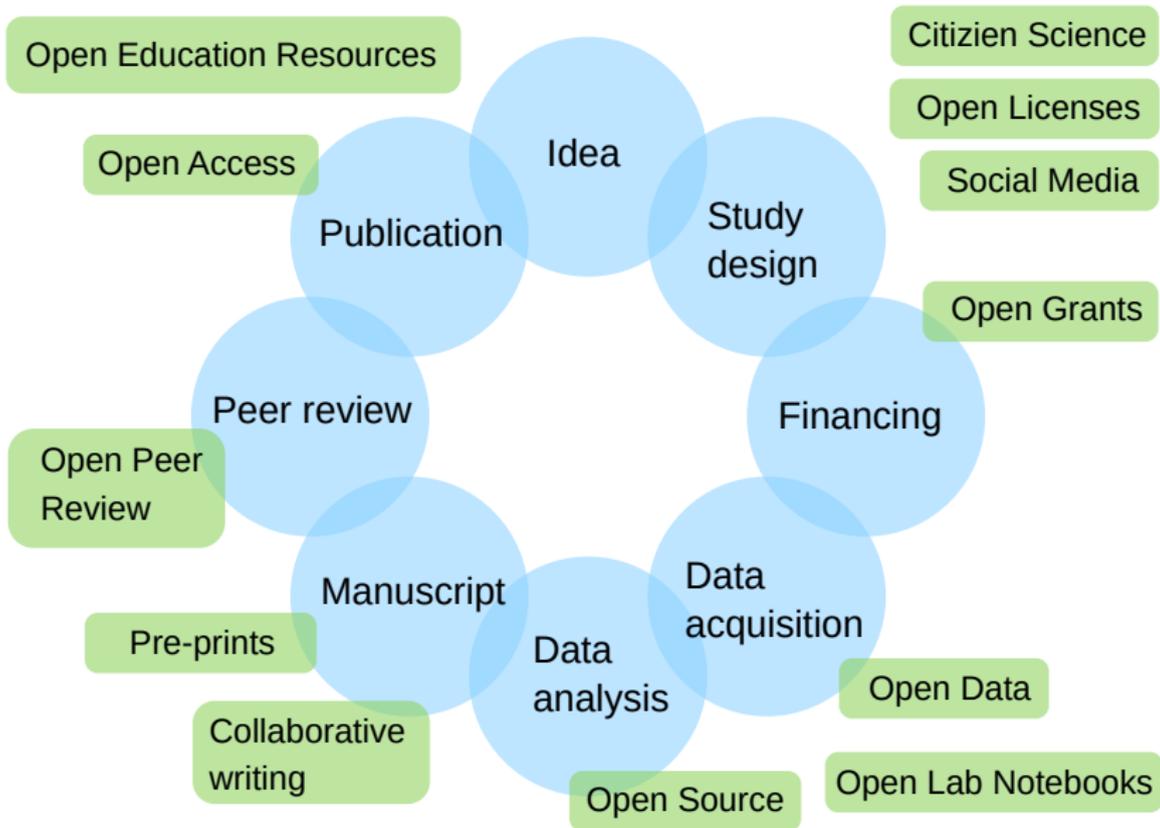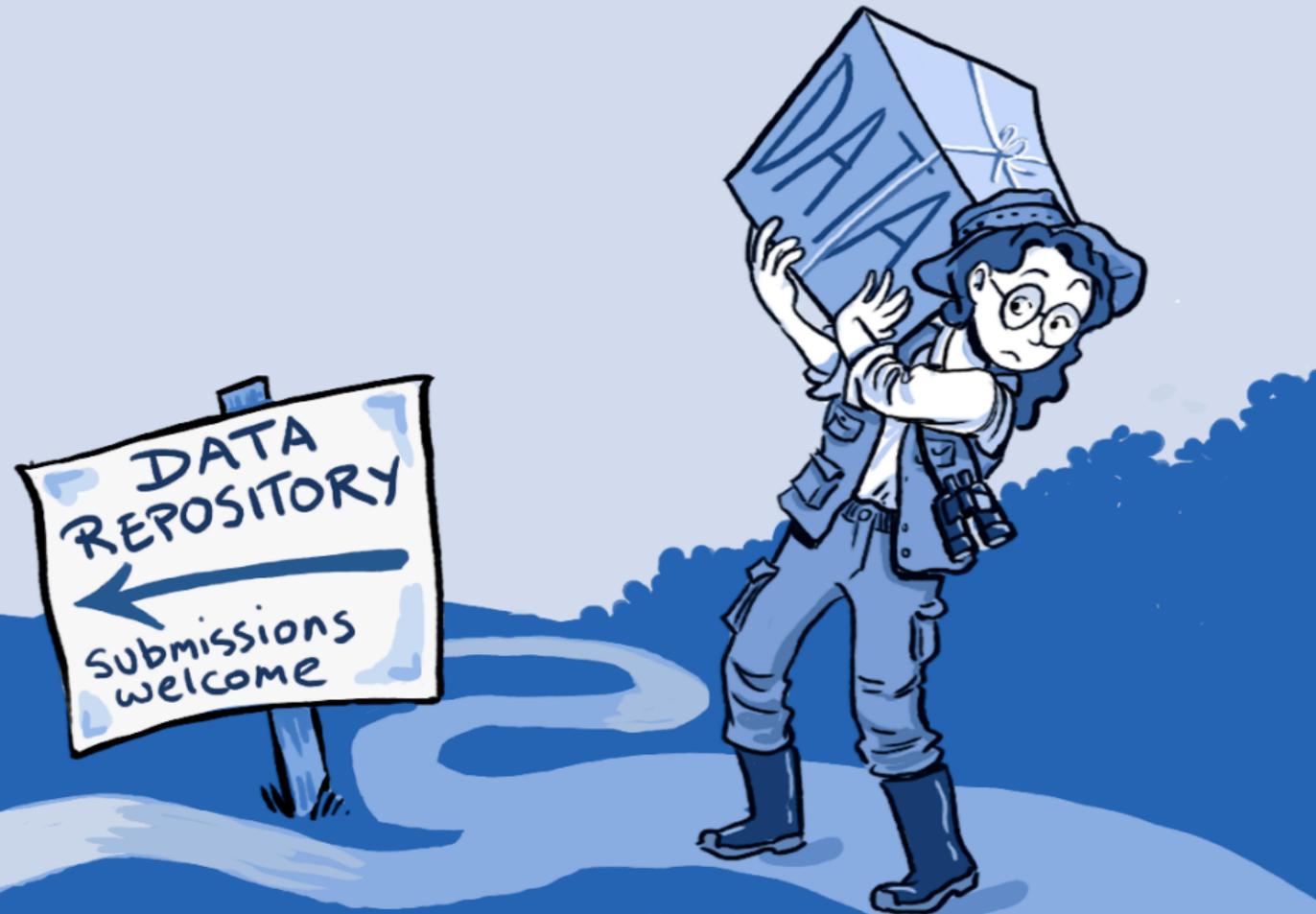
Konrad U. Förstner

ZB MED – Information Center Life Sciences, Cologne, Germany &
TH Köln, Cologne Germany

Novel approaches in Open Science
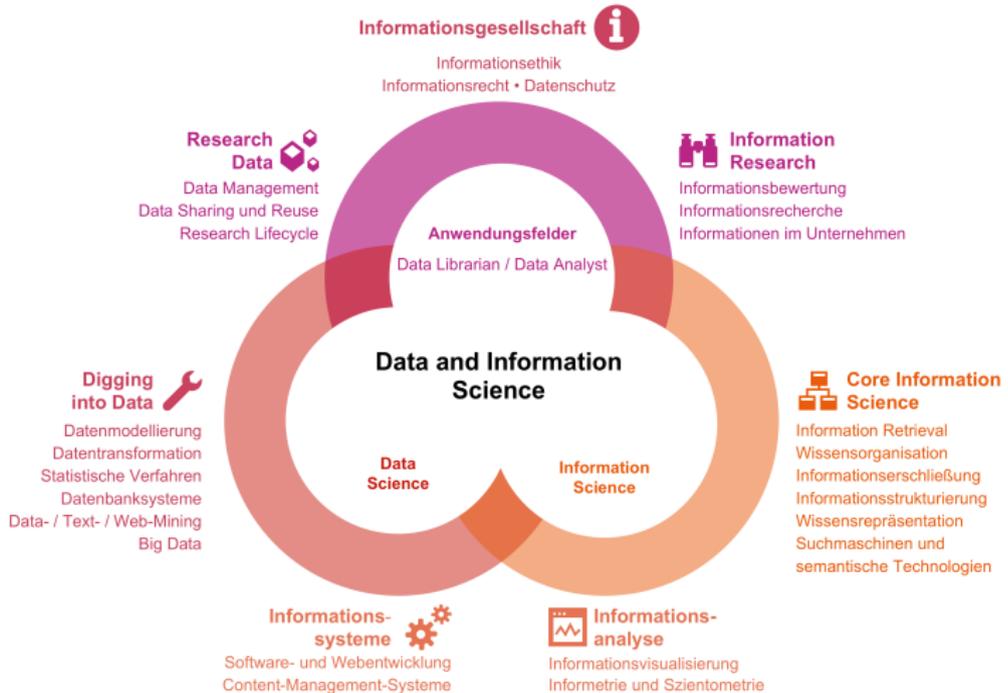University Library, University of Mannheim
Oct. 22$^{nd}$ 2018

Idea

Study design

Financing

Data acquisition

Data analysis

Manuscript

Peer review

Publication

Open Education Resources

Open Access

Citizien Science

Open Licenses

Social Media

Open Grants

Open Peer Review

Pre-prints

Collaborative writing

Open Source

Open Data

Open Lab Notebooks

illustration Ainsley Seago; https://doi.org/10.1371/journal.pbio.1001779

# Shifting to Data Savvy:
# The Future of Data Science
# In Libraries

# New Bachelor at TH Köln: *Data and Information Science*



**Informationsgesellschaft**
Informationsethik
Informationsrecht • Datenschutz

**Research Data**
Data Management
Data Sharing und Reuse
Research Lifecycle

**Information Research**
Informationsbewertung
Informationsrecherche
Informationen im Unternehmen

**Anwendungsfelder**
Data Librarian / Data Analyst

**Data and Information Science**

**Digging into Data**
Datenmodellierung
Datentransformation
Statistische Verfahren
Datenbanksysteme
Data- / Text- / Web-Mining
Big Data

**Core Information Science**
Information Retrieval
Wissensorganisation
Informationserschließung
Informationsstrukturierung
Wissensrepräsentation
Suchmaschinen und
semantische Technologien

**Data Science**

**Information Science**

**Informations-systeme**
Software- und Webentwicklung
Content-Management-Systeme

**Informations-analyse**
Informationsvisualisierung
Informetrie und Szientometrie

Currently under contruction at the ZBIW of the TH Köln:

Certificate course *Data Librarian*

*Software and data skills for library professionals*

# WELCOME TO LIBRARY CARPENTRY

LEARN MORE

## WHAT WE DO

Lessons

Workshops

Community

**Data intro for librarians**

*An introduction to data structures, regular expressions, and computing terms*

**Unix Shell**

*An introduction to command line interfaces and task automation using the Unix shell*

**OpenRefine**

*An introduction to cleaning up and enhancing a dataset using OpenRefine*

**Git Intro for Librarians**

*An introduction to version control using Git and GitHub for collaboration*

**SQL for Librarians**

*An introduction to relational database management using the SQLite tool*

**Webscraping**

*An introduction to extracting structured data from websites using a range of tools*

**Tidy data for librarians**

*An introduction to good data organisation, which is the foundation of much of our day-to-day work in libraries.*

**Introduction to Python**

*An introduction to Python, a general purpose programming language*

**Data Intro for Archivists**

*An introduction to data structures, regular expressions, and computing terms for archivists*

Science ⇌ Technology

Software
- an ubiquitous research tool

It is unquestionable that there is a strong and growing dependence of research on software.

Software is also a result of the scientific work.

Quality, accessibility, citability, etc. have to be ensured.

The importance of software for research is widely ignored.

# Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates.

Eklund A[1], Nichols TE[2], Knutsson H[3].

⊕ **Author information**

**Erratum in**

Correction for Eklund et al., Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates. [Proc Natl Acad Sci U S A. 2016]

## Abstract

The most widely used task functional magnetic resonance imaging (fMRI) analyses use parametric statistical methods that depend on a variety of assumptions. In this work, we use real resting-state data and a total of 3 million random task group analyses to compute empirical familywise error rates for the fMRI software packages SPM, FSL, and AFNI, as well as a nonparametric permutation method. For a nominal familywise error rate of 5%, the parametric statistical methods are shown to be conservative for voxelwise inference and invalid for clusterwise inference. Our results suggest that the principal cause of the invalid cluster inferences is spatial autocorrelation functions that do not follow the assumed Gaussian shape. By comparison, the nonparametric permutation test is found to produce nominal results for voxelwise as well as clusterwise inference. These findings speak to the need of validating the statistical methods being used in the field of neuroimaging.

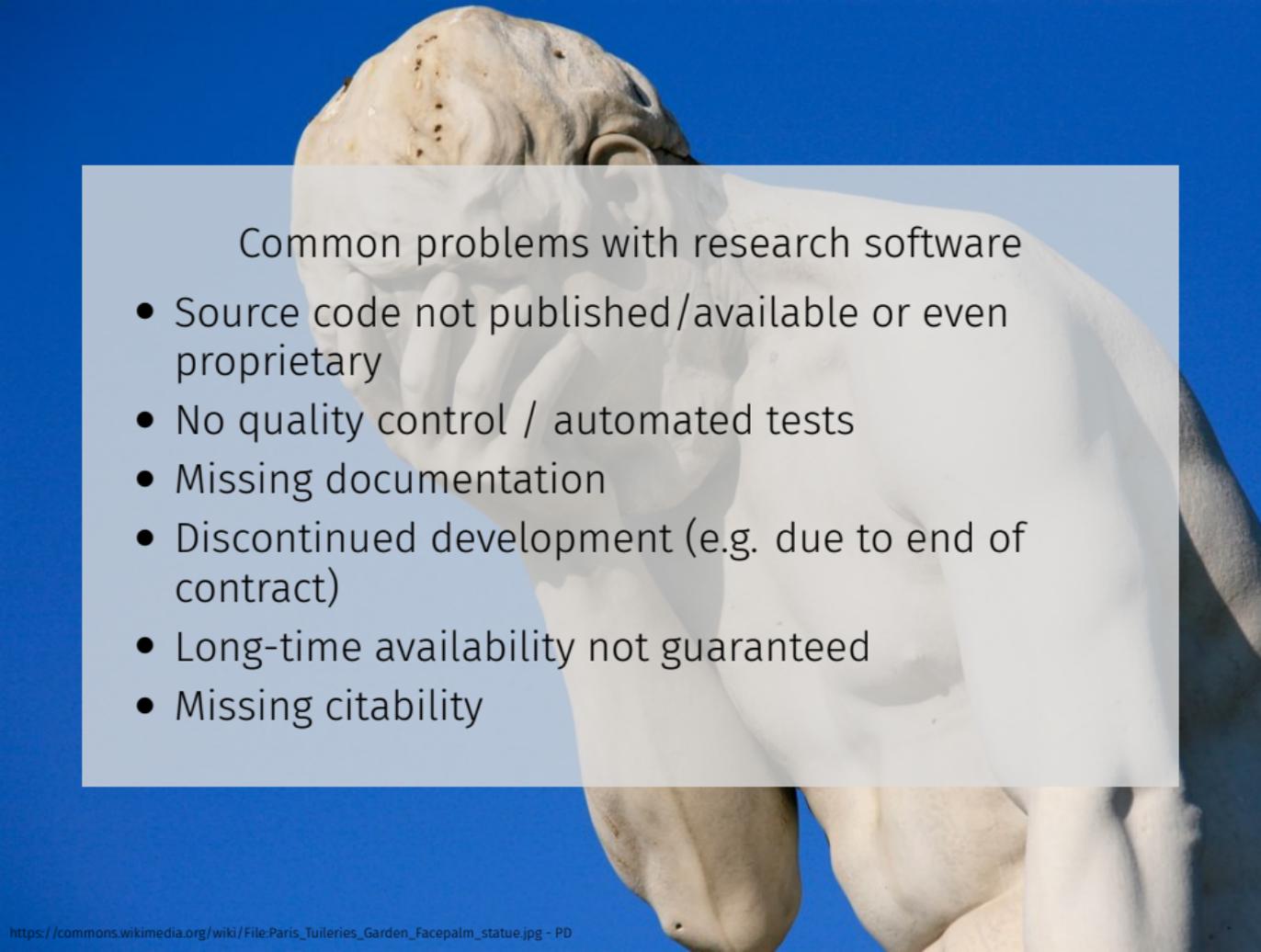# Does high public debt consistently stifle economic growth? A critique of Reinhart and Rogoff

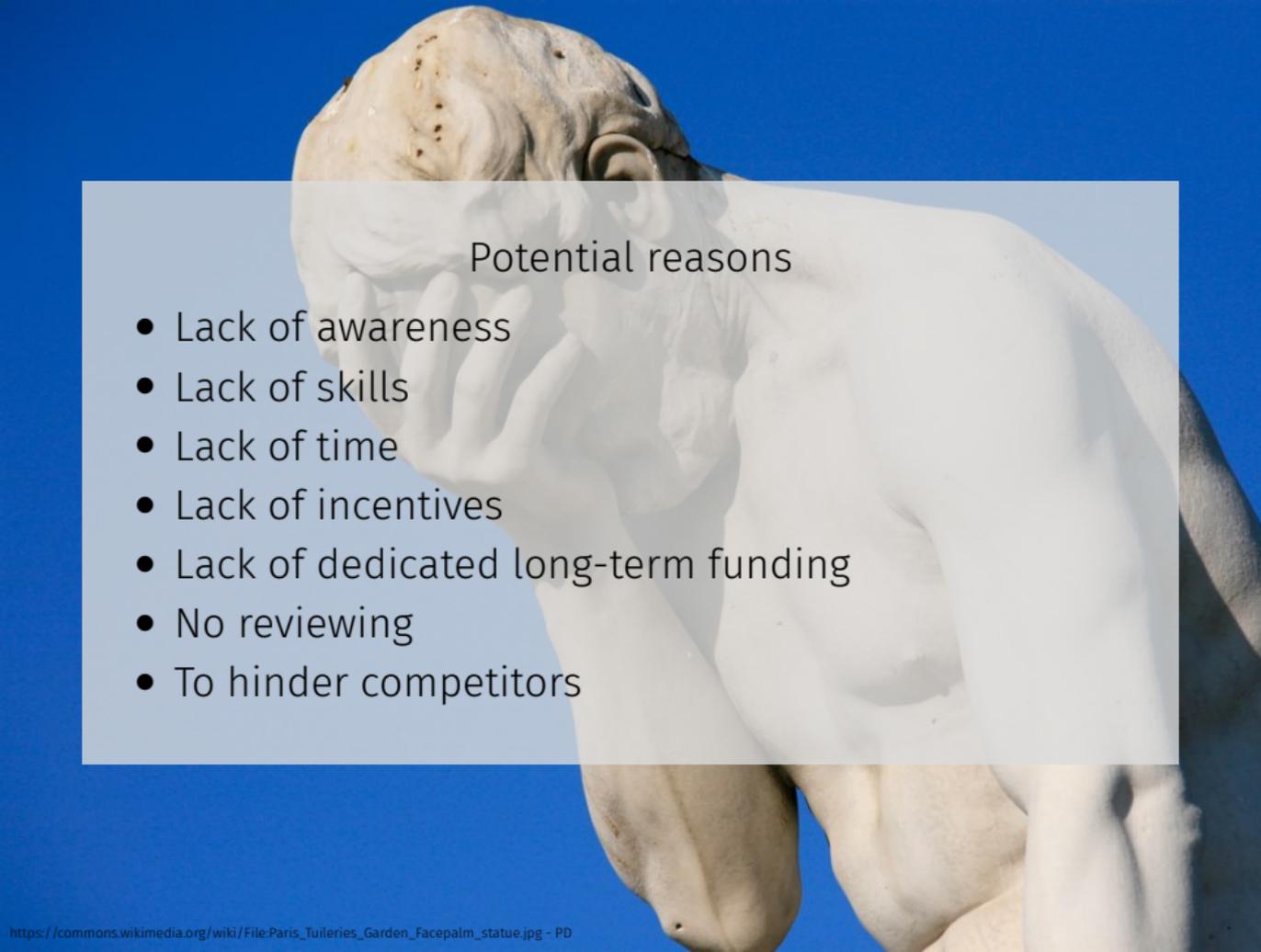Thomas Herndon, Michael Ash, Robert Pollin

## Abstract

We replicate Reinhart and Rogoff (2010A and 2010B) and find that selective exclusion of available data, coding errors and inappropriate weighting of summary statistics lead to serious miscalculations that inaccurately represent the relationship between public debt and GDP growth among 20 advanced economies. Over 1946–2009, countries with public debt/GDP ratios above 90% averaged 2.2% real annual GDP growth, not -0.1% as published. The published results for (i) median GDP growth rates for the 1946–2009 period and (ii) mean and median GDP growth figures over 1790–2009 are all distorted by similar methodological errors, although the magnitudes of the distortions are somewhat smaller than with the mean figures for 1946–2009. Contrary to Reinhart and Rogoff's broader contentions, both mean and median GDP growth when public debt levels exceed 90% of GDP are not dramatically different from when the public debt/GDP ratios are lower. The relationship between public debt and GDP growth varies significantly by period and country. Our overall evidence refutes RR's claim that public debt/GDP ratios above 90% consistently reduce a country's GDP growth.

Common problems with research software

- Source code not published/available or even proprietary
- No quality control / automated tests
- Missing documentation
- Discontinued development (e.g. due to end of contract)
- Long-time availability not guaranteed
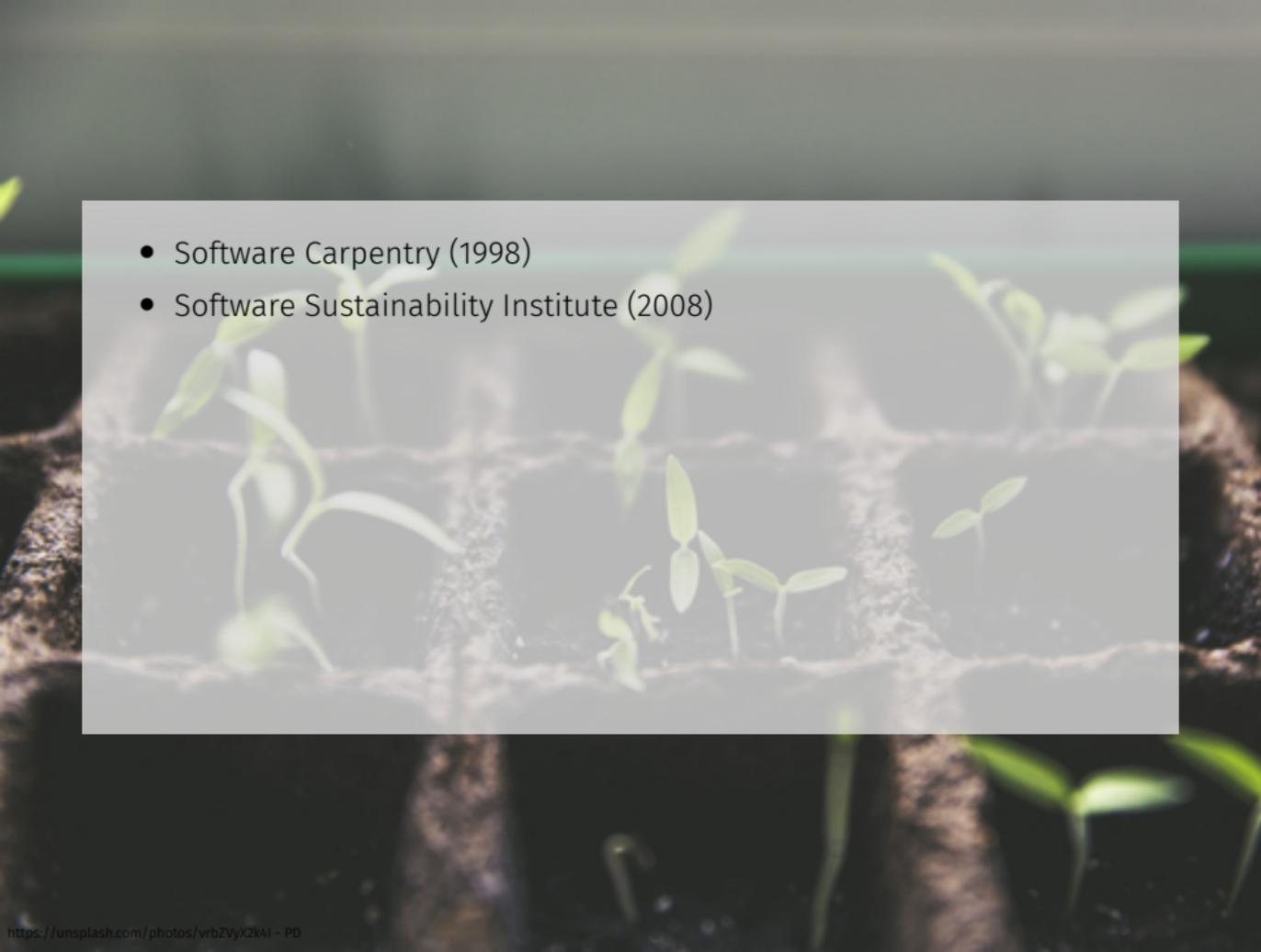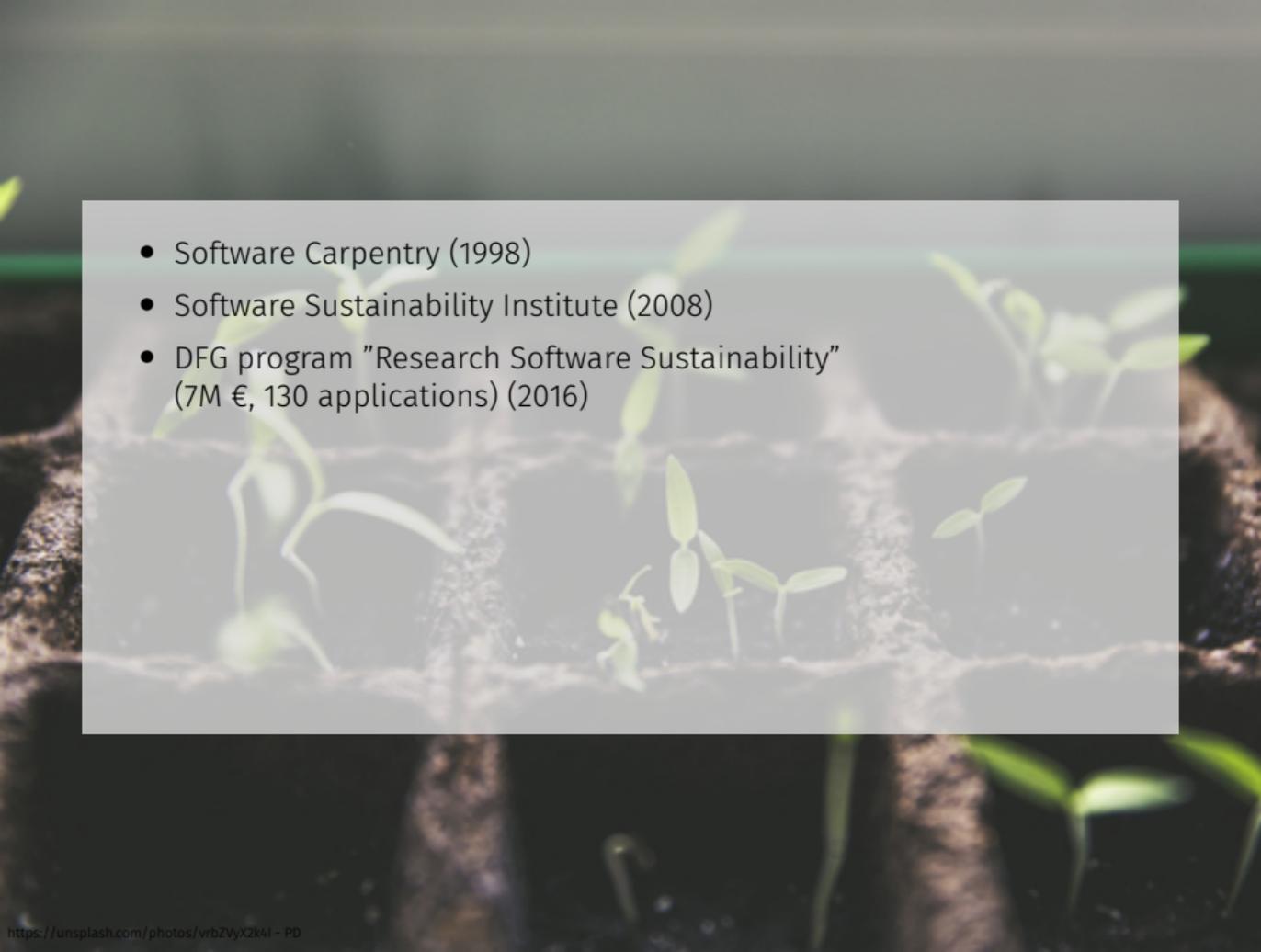- Missing citability

Potential reasons

- Lack of awareness
- Lack of skills
- Lack of time
- Lack of incentives
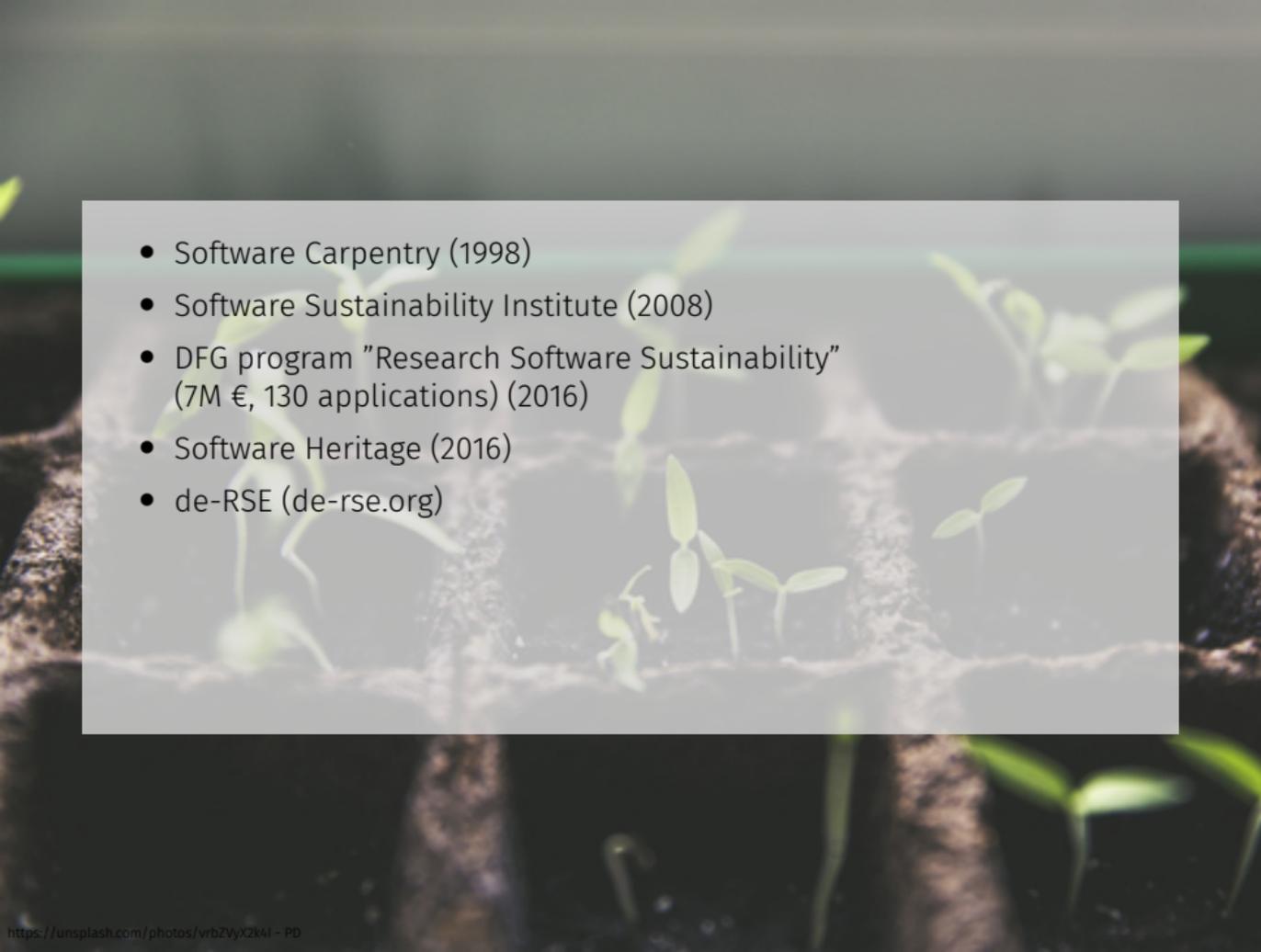- Lack of dedicated long-term funding
- No reviewing
- To hinder competitors

Several iniatives have been launched to address these issues.

- Software Carpentry (1998)

- Software Carpentry (1998)
- Software Sustainability Institute (2008)

- Software Carpentry (1998)
- Software Sustainability Institute (2008)
- DFG program "Research Software Sustainability" (7M €, 130 applications) (2016)

- Software Carpentry (1998)
- Software Sustainability Institute (2008)
- DFG program "Research Software Sustainability" (7M €, 130 applications) (2016)
- Software Heritage (2016)
- de-RSE (de-rse.org)

- Software Carpentry (1998)
- Software Sustainability Institute (2008)
- DFG program "Research Software Sustainability" (7M €, 130 applications) (2016)
- Software Heritage (2016)
- de-RSE (de-rse.org)
- Working group "Digital tool - Software and Service" as part of the focus initiative "Digital Information" of the Alliance of Science Organizations in Germany
- Several more ...

Recommendations on the

# Development, Use and Provision of Research Software

Research Software Working Group
in the Priority Initiative Digital Information
of the Alliance of German Science Organisations

Guiding principle

The concept of Good Scientific Practice (GSP) must be also applied to research software.

But what can Good Scientific Practice mean
for research software?

FAIR principle should also be applied to software

- Findable
- Accessible
- Interoperable
- Re-usable

Open

- Proper OSI conform licence

Three types of software

1. Small tools written for single purpose
2. Software applications (as research output)
3. Infrastructure and online services

All three levels are relevant and
have to be addressed.

Exact needs and possibilities might differ between scientific communities.

Discourse must also happen inside of these communities.

E.g. what exactly means "reproducibility" (bit-identical compilation?) and how long would this needed to be guaranteed?

Raise the awareness for the relevance of research software.

Introduce standards and mechanisms for quality control of research software.

Create institutional platforms to publish and archive software/source code/workflows.

Enable citation of such items.
e.g. using Citation File Format
(CFF - citation-file-format.github.io)

Make these citations part of the scientific reputation system.

Foster the education of computational skills inside of the scientific community.

Develop new carreer paths like
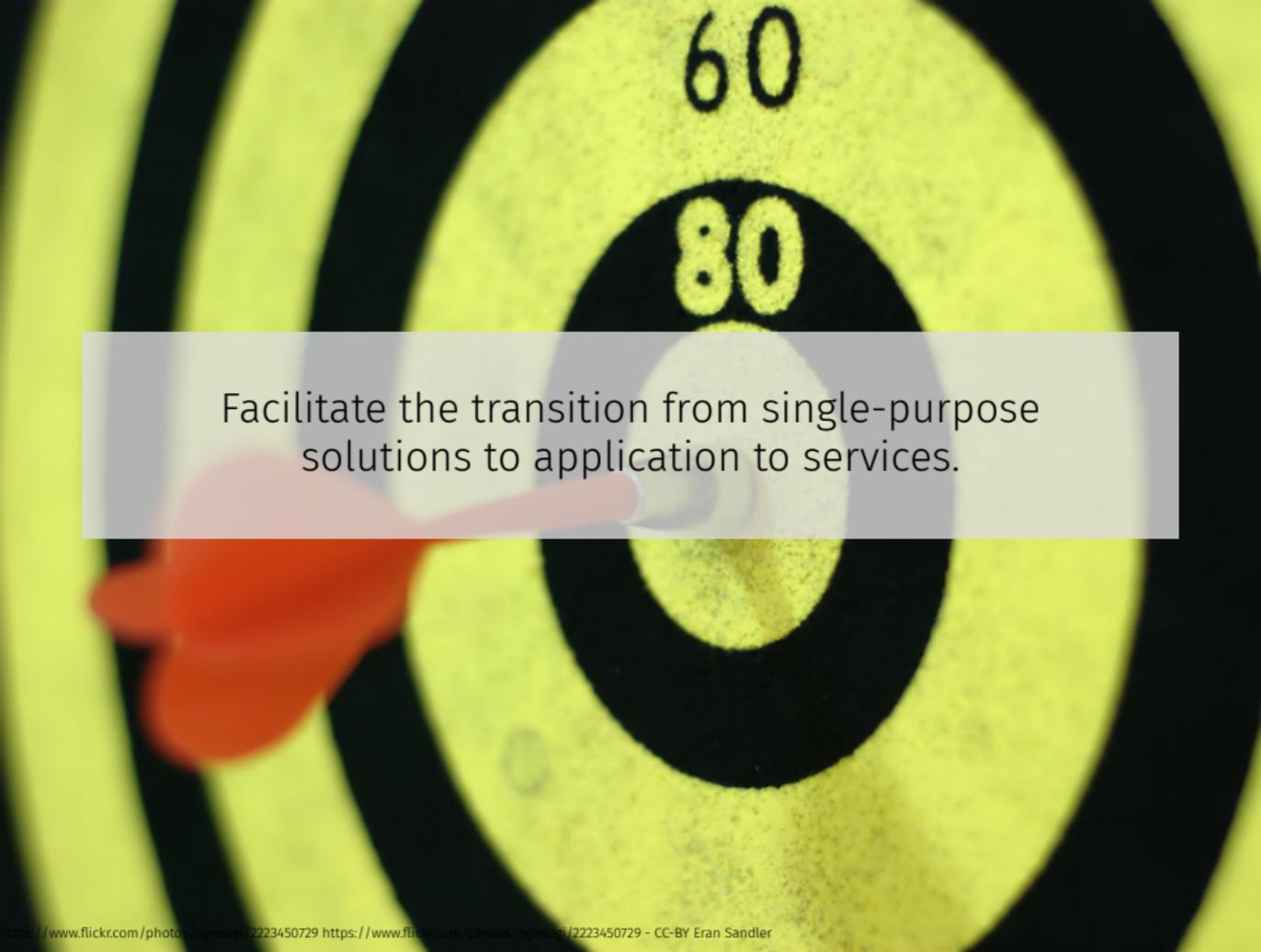Research Software Engineers,
Software Librarians, Data Scientists.

Raise awareness about and teach legal aspects (i.e. licensing) of software.

Raise awareness about and teach legal aspects (i.e. licensing) of software.

Make open source the default.

Facilitate the transition from single-purpose solutions to application to services.

Provide long-term funding to enable sustainable software development.

A lot to do and a lot of open questions.

# www.allianzinitiative.de

Matthias Katerbow - Deutsche Forschungsgemeinschaft
Michael Goedicke - Deutsche Forschungsgemeinschaft
Leander Seige - Deutsche Forschungsgemeinschaft
Zeki Mustafa Dogan - Deutsche Forschungsgemeinschaft
Dirk Eisengräber-Pabst - Fraunhofer-Gesellschaft
Uwe Konrad - Helmholtz-Gemeinschaft
Bernadette Fritzsch - Helmholtz-Gemeinschaft
Björn Brembs - Hochschulrektorenkonferenz
Klaus Wannemacher - Hochschulrektorenkonferenz
Thomas Dandekar - Hochschulrektorenkonferenz
Georg Feulner - Leibniz-Gemeinschaft
Jürgen Fuhrmann - Leibniz-Gemeinschaft
Michael Franke - Max-Planck-Gesellschaft
Stefan Janosch - Max-Planck-Gesellschaft
Johannes Reetz -Max-Planck-Gesellschaft
Thomas Rode - Leopoldina
Mathias Bornschein - Bibliotheken der Ressortforschungseinrichtungen des Bundes

ZB MED

TH Köln

@konradfoerstner