

Software als Forschungsdaten

aus der Perspektive der Software-Technik-Forschung

Forschungsdaten-Soiree #2

Anne Koziolk

MODELLING FOR CONTINUOUS SOFTWARE ENGINEERING,
INSTITUT FÜR INFORMATIONSSICHERHEIT UND VERLÄSSLICHKEIT, KIT-FAKULTÄT FÜR INFORMATIK



Diese Folien sind auf <https://are.ipd.kit.edu/people/anne-koziolk/> zu finden

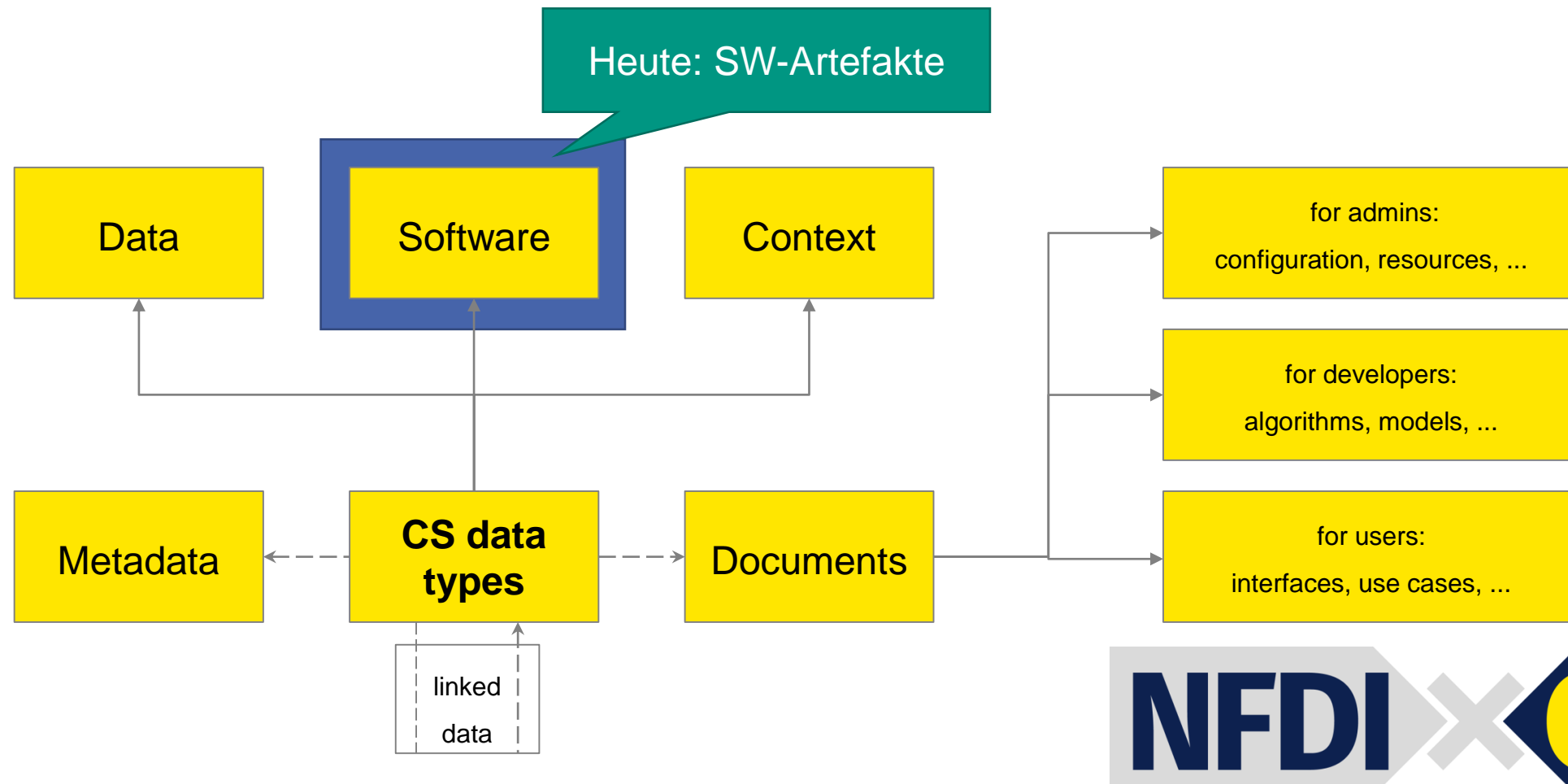
Überblick

- Welche Forschungsdaten?
- Wo stehen wir heute?
 - NISO Badges
 - Einschätzung zum Stand der FAIR Prinzipien und Herausforderungen
- Lösungsansätze in NFDIxCS

Dabei meine Perspektive: Software-Technik-Forschung

Für die Diskussion später: Passt es zu anderen Gebieten,
in denen Software Forschungsdaten sind?

Forschungsdaten in der Informatik



Welche Forschungsdaten betrachten wir?

■ Von Forschern erzeugt

- Software
 - Skripte für Datenanalyse
 - Programme, die neue Algorithmen umsetzen
 - Software-Tools
- Modelle
- Andere Daten

■ Beobachtungen

- Unter Beteiligung von Menschen oder Umgebung beobachtet
- Automatisiert erzeugt

■ Von anderen erzeugte Eingabedaten: u.a. auch Software



Fokus im Folgenden:
SW-Artefakte

Artefakte [ACM Badging,
<https://www.acm.org/publications/policies/artifact-review-and-badging-current>]

Wo stehen wir heute in der SWT-Forschung?

- Zunehmende Bedeutung
- „Artifact Evaluation Track“ auf den meisten SWT-Konferenzen
- ACM Artifact Reviewing and Badging
 - “An experimental result is not fully established unless it can be independently reproduced.”
 - “By "artifact" we mean a digital object that was either created by the authors to be used as part of the study or generated by the experiment itself. For example, artifacts can be software systems, scripts used to run experiments, input datasets, raw data collected in the experiment, or scripts used to analyze results.”
- Aber
 - Veröffentlichung der Artefakte keine Voraussetzung für Publikationen
 - Wiederverwendung von Artefakten weiterhin schwierig

NISO Artifact Reviewing and Badging

Open Research Object



Author-created digital objects used in the research (including data and code) are permanently archived in a public repository that assigns a global identifier and guarantees persistence, and are made available via standard open licenses that maximize artifact availability.

Research Object Reviewed (ROR)



All relevant author-created digital objects used in the research (including data and code) were reviewed according to the criteria provided by the badge issuer.

The badge metadata should link to the award criteria.

(possibly different levels)

Results Reproduced (ROR-R)



An additional step was taken or facilitated by the badge issuer (e.g., publisher, trusted third-party certifier) to certify that an independent party has regenerated computational results using the author-created research objects, methods, code, and conditions of analysis.

Results Replicated (RER)



An independent study, aimed at answering the same scientific question, has obtained consistent results leading to the same findings (potentially using new artifacts or methods).

The badge links to the persistent identifier for that secondary publication. This badge is awarded by the publisher of the original work that is being badged.

Stand und Herausforderungen FAIRer Forschungsdaten in SWT

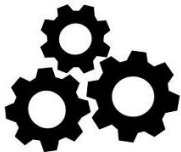
F
indable



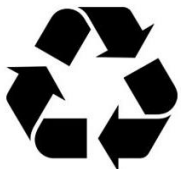
A
ccessible



I
nteroperable



R
eusable



- DOIs oder Links in Papers genannt.
- Archivierung oft Github, Zenodo o.ä.
- **Aber** wenig Nutzung von Metadaten, keine Suchmaschinen
- Open Source etabliert
- **Aber** Anreize fehlen, um Aufwand auszugleichen.
- Was bedeutet das für Artefakte aus der Software-Entwicklung?
- Ausführbarkeit/Portabilität:
 - Verfügbare Technologien: Container, virtuelle Maschinen
 - **Aber**: Ganze VMs zu archivieren ist nicht die Lösung, wie mit Abhängigkeiten umgehen?
- Prinzipielle Techniken gerade in SWT bekannt
- **Aber** SW-Artefakte trotzdem schwer zu verstehen und anzupassen

Image „[FAIR data principles](#)“ by [SangyaPundir](#) is licensed under [CC-BY-SA-4.0](#) on Wikipedia
s. auch Hasselbring et al., „From FAIR research data toward FAIR and open research software“, DOI: 10.1515/itit-2019-0040

(Einige) Ansätze in NFDIxCS für SW

- Findable
 - Establish version-aware metadata and semantics beyond paper PDFs
 - Accessible
 - Legal support for open source licensing
 - Citation and attribution: Include data publications in DBLP
 - Publication processes: Integrate in publication system
 - Interoperable
 - Long term archiving including dependencies
 - Reusable execution environments
 - Reusable
 - Review system for assessing quality
- Bereitstellung und Nutzung von SW als Forschungsdaten vereinfachen

Überblick

- Welche Forschungsdaten?
- Wo stehen wir heute?
 - NISO Badges
 - Einschätzung zum Stand der FAIR Prinzipien und Herausforderungen
- Lösungsansätze in NFDIxCS
 - Ziele: Bereitstellung und Nutzung von SW als Forschungsdaten vereinfachen

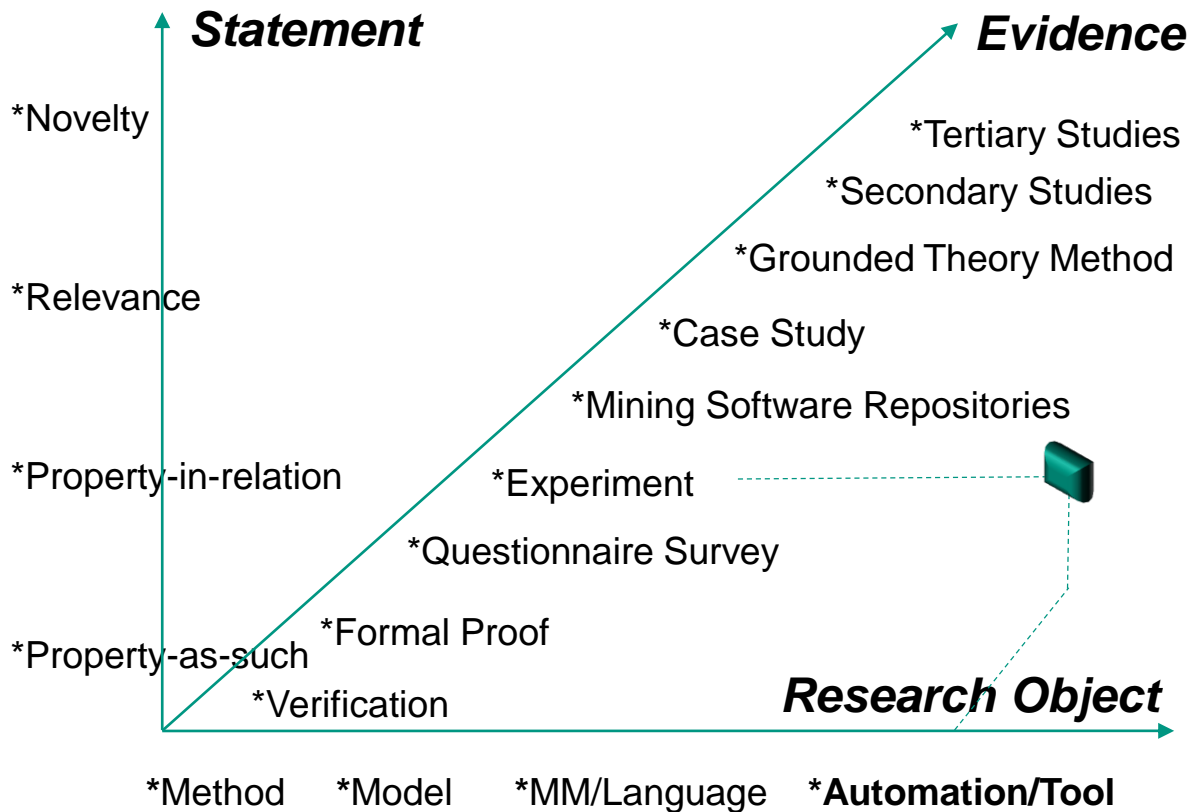
Dabei meine Perspektive: Software-Technik-Forschung

Für die Diskussion später: Passt es zu anderen Gebieten, in denen Software Forschungsdaten sind?

Diese Folien sind auf <https://are.ipd.kit.edu/people/anne-koziolk/> zu finden

- [ACM Badging] „By "artifact" we mean a digital object that was either created by the authors to be used as part of the study or generated by the experiment itself. For example, artifacts can be software systems, scripts used to run experiments, input datasets, raw data collected in the experiment, or scripts used to analyze results.” <https://www.acm.org/publications/policies/artifact-review-and-badging-current>

Aktuelle Forschung: Validitäts-Aussagen in der SWT-Forschung



Data Item	Data Description
Research Object	Investigated object of research (dt. „Untersuchungsgegenstand“)
Statement	Kind of statement: (i) One paper can be related to many statements. (ii) Each statement may have given different kinds and strength of evidence. (iii) Papers may add evidence to a statement of other papers with additional data or arguments. (iv) Papers may contain statements that refine statements of other papers or (v) even contradict them with some evidence.
Evidence	Research method design and evidence of statement validity
Validation Question (VQ)	Template for validation question: * e.g., „How <i>effective</i> is Tool A in comparison to Tool B?“

Gemeinsame Arbeit mit Angelika Kaplan und Ralf Reussner, Folie von Angelika Kaplan

ACM Artifact Reviewing and Badging

Functional



Artifacts documented, consistent, complete, exercisable, and include appropriate evidence of verification and validation

Reusable



Functional + very carefully documented and well-structured to the extent that reuse and repurposing is facilitated. In particular, norms and standards of the research community for artifacts of this type are strictly adhered to.

Available



Placed on a publically accessible archival repository. A DOI or link to this repository along with a unique identifier for the object is provided.

Replicated



Main results of the paper have been independently obtained in a subsequent study by a person or team other than the authors, without the use of author-supplied artifacts.

Reproduced



Main results of the paper have been obtained in a subsequent study by a person or team other than the authors, using, in part, artifacts provided by the author.

<https://www.acm.org/publications/policies/artifact-review-badging>

SWT Adaption der ACM Badges

Functional

No Badge

Reusable



Available



Replicated



Reproduced



Artifacts documented, consistent, complete, exercisable, and include appropriate evidence of verification and validation

Functional + very carefully documented and well-structured to the extent that reuse and repurposing is facilitated. In particular, norms and standards of the research community for artifacts of this type are strictly adhered to.

Functional + placed on a publicly accessible archival repository. A DOI or link to this repository along with a unique identifier for the object is provided.

Available + main results of the paper have been obtained in a subsequent study by a person or team other than the authors, using, *in part*, artifacts provided by the author.

Available + the main results of the paper have been independently obtained in a subsequent study by a person or team other than the authors, *without the use of author-supplied artifacts*.

vom Artifact Evaluation Track der ICSE Konferenz, <https://conf.researchr.org/track/icse-2020/icse-2020-Artifact-Evaluation#Call-for-Submissions>