# (Re)Use of Research Results … why should we?

Maria Teresa Baldassarre

Department of Informatics – University of Bari

mariateresa.baldassarre@uniba.it

*Software Engineering Research LABoratory*

# Who am I



✉ mariateresa.baldassarre@uniba.it

🐦 @mtbaldassarre

⇨ **Associate Professor** at the Department of Informatics - University of Bari (www.di.uniba.it)

⇨ **Coordinator of the «Process & Product Quality»** area @Software Engineering Research LAB (serlab.di.uniba.it)

⇨ **Quality Manager** @SER&Practices Spin-Off (https://serandp.com/en/)

⇨ **Member** of the International Software Engineering Reserch Netwok (**ISERN**)

# CS Department –
# Dipartimento di Informatica
# BARI - Puglia

⇨ Main research interests:

- ❑ SOFTWARE PROCESS AND PRODUCT QUALITY
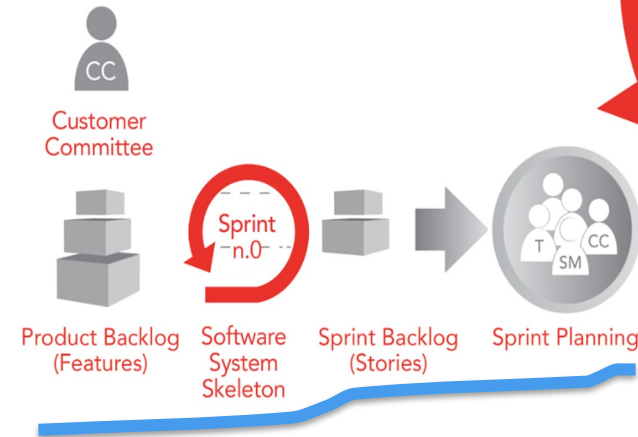- ❑ HUMAN FACTORS IN SOFTWARE ENGINEERING
- ❑ EMPIRICAL SOFTWARE ENGINEERING

Software Engineering Research LABoratory

# Agile User&Quality Oriented Development...

# …Agile User&Quality Oriented Development

- **BAMBOO**: continuous integration/deployment
- **JIRA**: application lifecycle management
- **SONARQUBE**: quality management
- **SVN**: version control system

**Software Engineering Research & Practices**

Spin-off of the University of Bari - established in 2006.

30 employees

- ❑ 9001:2008 - Quality management systems - Requirements
- ❑ 14001: 2004 - Environmental management systems
- ❑ 25000:2014 - Systems and software engineering – **First in Italy to assess certification of a sofware product**

9001:2008

14001:2004

25000:2014

- ❑ SERLab carries out research and empirical validation of results
- ❑ SER&P transfers the results of these activities to industry; provides data and industrial context for field experimentation



SOFTWARE SYSTEM DESIGN & DEVELOPMENT

SOFTWARE SYSTEM GOVERNANCE & SECURITY

SOFTWARE PROCESS & PRODUCT QUALITY

PROJECT MANAGEMENT

# RESEARCH COLLABORATIONS

# INDUSTRIAL COLLABORATIONS

# Is it important for a scientist to Report Research Results so others can (Re)Use them?

*" … the ideas we can most trust are those that have been the most tried and tested.*

*For that reason many of us are involved in this process called 'science' which produces trusted knowledge by sharing one's ideas and trying out and testing the ideas of others … "*

cit. Popper

# Produce & Report research results

# ReUse results/findings …



… to improve reproducibility and transparency

# «RESULTS PARADOX»

# «RESULTS PARADOX»

**«FACTS & TRUTH»**

Keep research results at arm's length

Objective investigator – detective

Follows data with discipline; never indulges in data massaging or cherry picking

**«BE PERSUASIVE»**

Pressure of publishing clear novel and positive findings on behalf of funding agencies, evaluation committees

Good lawyer

Arguments and produces amounts of beautiful and convincing results

Chambers, C.D., Tzavella, L. The past, present and future of Registered Reports. *Nat Hum Behav* **6**, 29–42 (2022).
https://doi.org/10.1038/s41562-021-01193-7

shutterstock.com · 155910122

⇨ Researchers attempt to solve this paradox … questionable research practices … reduce confidence of conclusions … harm reproducibility …

# Questionable Research Practices (QPRs) Hurt Science ...

### HARKing (Hypotheszing After Results are Known)

Neat data, what explains it?

- Acceptable in explanatory not confirmatory

### Post-hoc Rationalizing

Story-telling to explain the data found in a study

- Acceptable in explanatory/inductive theory building not confirmatory

John LK, Loewenstein G, Prelec D (2012) Measuring the prevalence of questionable research practices with incentives for truth telling. Psychol Sci 23(5):524–532. https://doi.org/10.1177/0956797611430953

# ... Questionable Research Practices Hurt Science

**File-drawer effect**

- Hmm, bad outcome, bin it. Negative result – reject. Not published. Do not appear in meta-analysis and SLRs

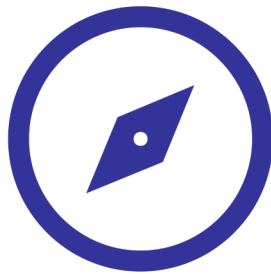**Forking paths in data analysis choices after seeing the data (Researcher Bias)**

- Let's use a Kruskal-Wallis test and then a Lewandoski-Neymar test of significance (instead of?)

QRPs result when publication **venue** and publication **significance/novelty** are emphasized over replication & soundness of the method

# Registered Reports

free researchers from the preasure to engage in QRPs

Avoid the RESULTS-ORIENTATION

Deal with RESEARCHER BIAS

Focus on SOUNDNESS OF THE RESEARCH PLAN & SIGNIFICANCE OF THE RESEARCH QUESTION

⇨ Ernst, N.A., Baldassarre, M.T. Registered reports in software engineering. *Empir Software Eng* **28**, 55 (2023). https://doi.org/10.1007/s10664-022-10277-5

# Registered Reports … why?

**Pre-registration (clinical trials)**: register your protocol including planned hypothesis, data collection, data analysis that is «registered» BEFORE the study is conducted

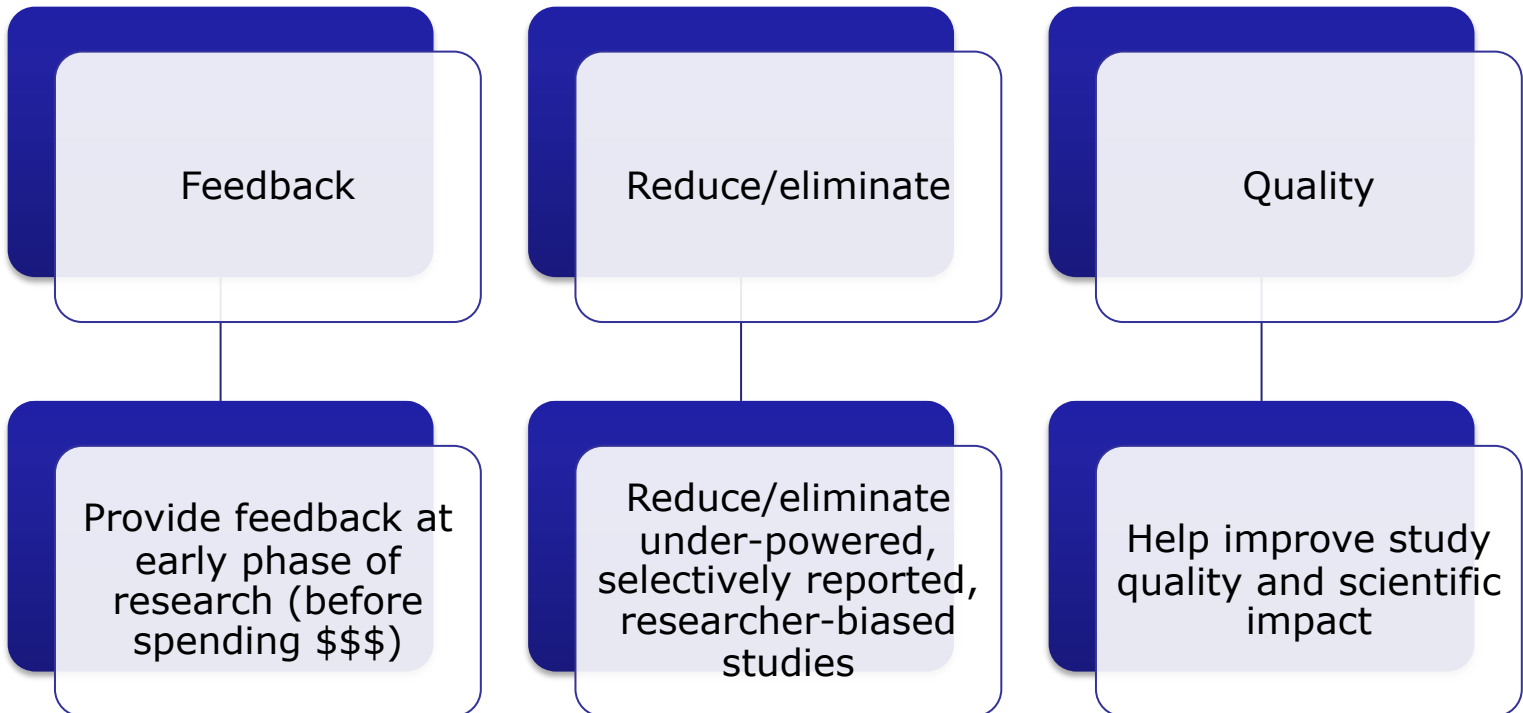**Protocol** comits to analysis and expected outcomes
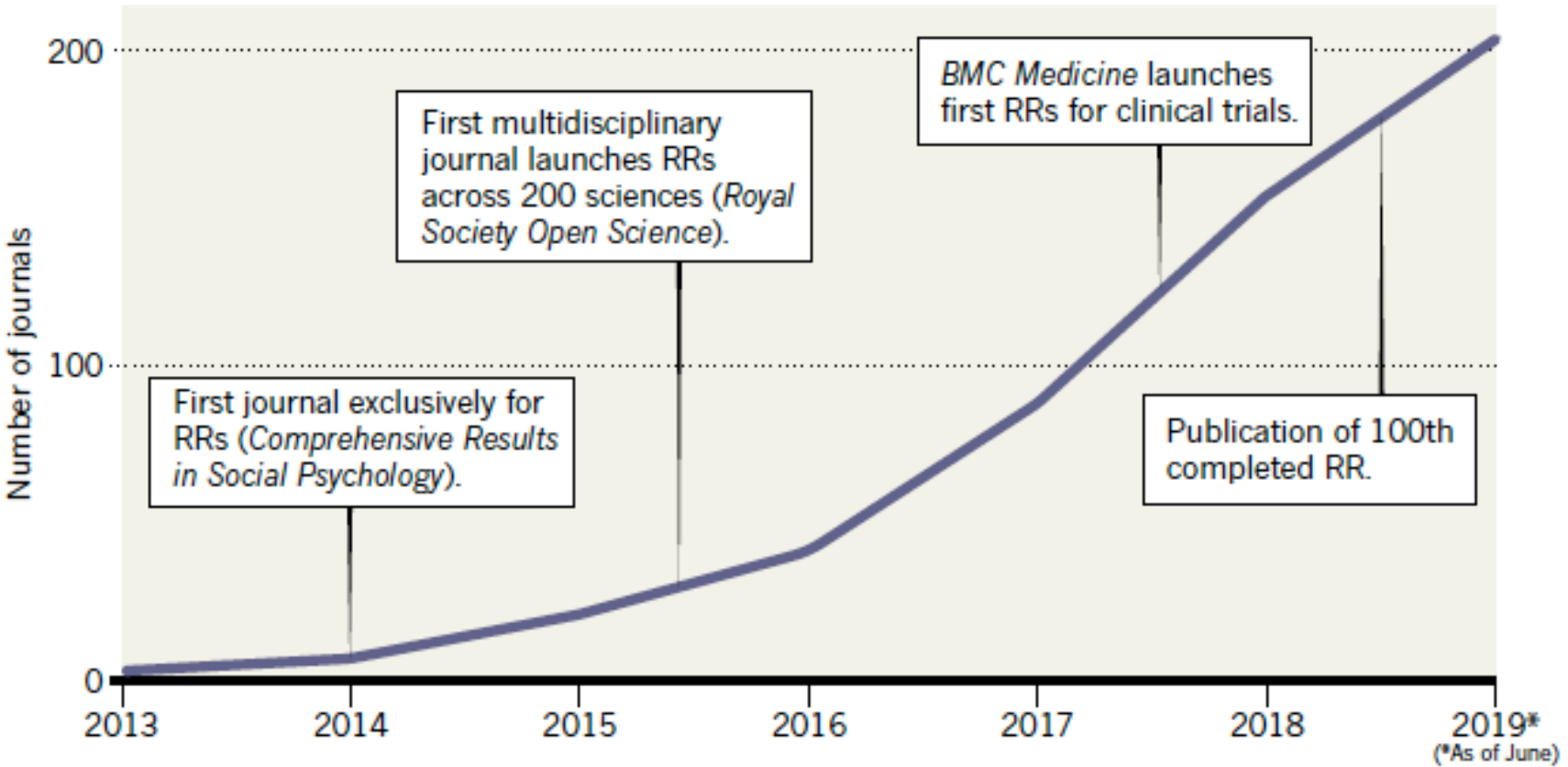
**Registered Report**: Peer-reviewed pre-registration

# … Registered Reports … why?

⇨ Benefits

| Feedback | Reduce/eliminate | Quality |
|---|---|---|
| Provide feedback at early phase of research (before spending $$$) | Reduce/eliminate under-powered, selectively reported, researcher-biased studies | Help improve study quality and scientific impact |

*Software Engineering Research LABoratory*

# RAPID RISE

Since 2013, the number of journals offering Registered Reports (RRs) has risen to more than 200 titles.



Number of journals

200

100

0

First journal exclusively for RRs (*Comprehensive Results in Social Psychology*).

First multidisciplinary journal launches RRs across 200 sciences (*Royal Society Open Science*).

BMC Medicine launches first RRs for clinical trials.

Publication of 100th completed RR.

2013    2014    2015    2016    2017    2018    2019*
(*As of June)

*Software Engineering Research LABoratory*

# RR in SW_Engineering

**EMSE J.** → MSR,
ICSME, then ESEM,
now CHASE, SANER,
ICPC
**TOSEM** (direct submit)
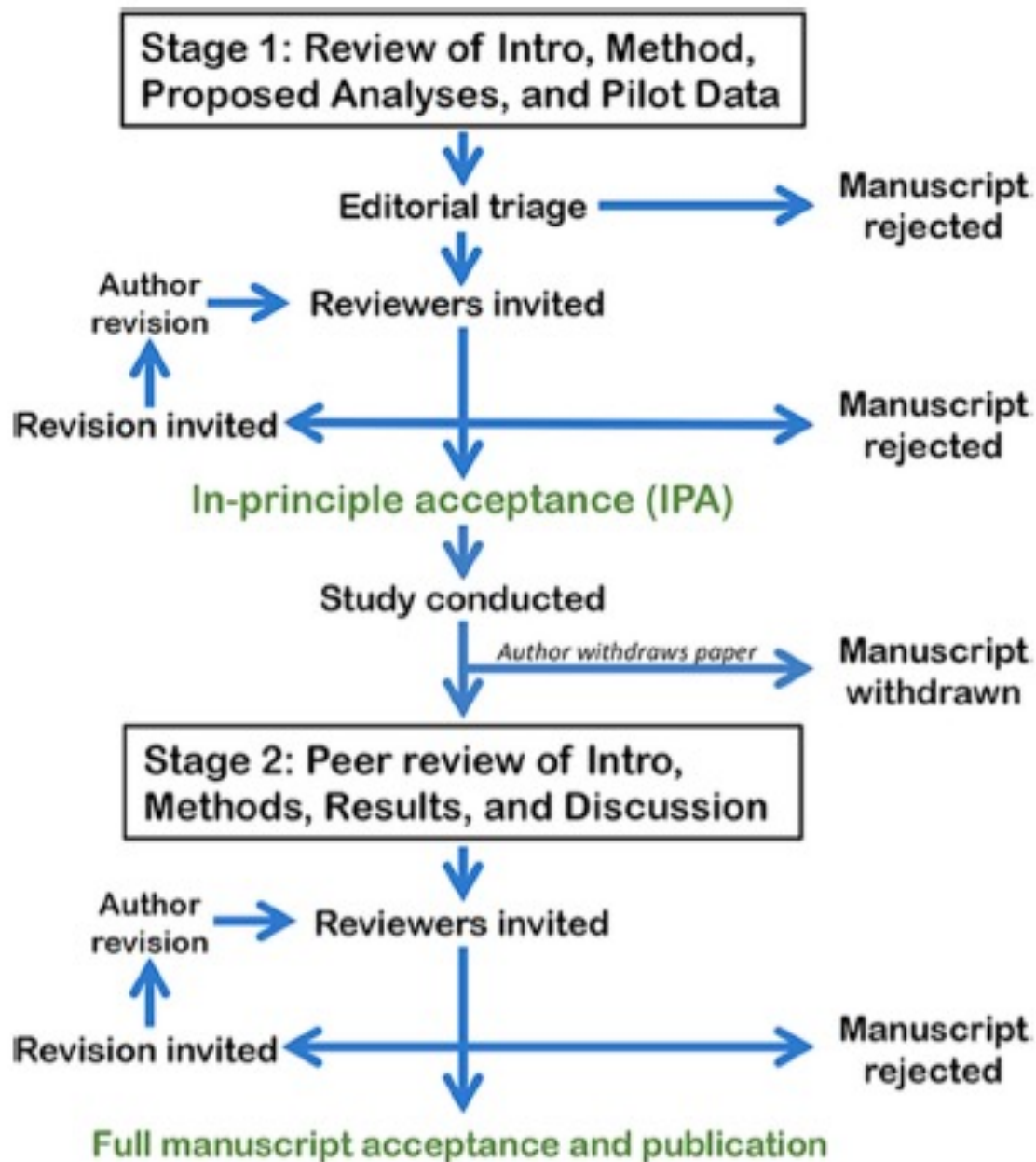**CSE** special issue

(ACM, Springer, T&F)

**Fig. 1** Stages of the Registered Reports workflow. Center for Open Science (https://www.cos.io/initiatives/registered-reports?#tabid3) CC-BY-NoDerivs 4.0

# Phase 1 – Review Criteria

**Is this study novel, significant, able to find effects?**

1. **Importance** of the research question(s).
2. Logic, rationale, and plausibility of the **proposed hypotheses**.
3. **Soundness** and **feasibility** of the methodology and analysis pipeline (including statistical power analysis where appropriate).
4. **Clarity** and degree of methodological detail for replication.
5. Will results obtained **test** the stated hypotheses?

# Phase 2 – Review Criteria

**Did the authors execute on Phase 1 plan?**

1. Whether the data are able to test the authors' proposed hypotheses by satisfying the approved outcome-neutral conditions (such as quality checks, positive controls)
2. Whether the Introduction, rationale and stated hypotheses **are the same** as the approved Stage 1 submission (required)
3. Whether the authors adhered **precisely** to the registered experimental procedures
4. Whether any unregistered post hoc analyses added by the authors are **justified**, methodologically sound, and informative
5. Whether the authors' conclusions are **justified** given the data

# Current state of RR in SE

MSR 2020 feedback on IPA:

"I think it is a key principle. However, in a way it also **raises the bar** significantly for the Registered Reports"

"[...] the fact that the results are missing, helps reviewers and authors **focus on the methodological issue**, which is a great added value in the review process [...]"

# MSR Results - IPA

"During my review, though, I had the feeling that **more interaction with the authors** could add even further value"

"I think the **EMSE paper still needs a careful assessment**, as it is still possible that the operation or the application of the protocol turns out to be wrong [...]"

"I felt a **bit uncomfortable to have this burden** on my shoulders as a reviewer so early in the process."

No (3 responses):

"A registered report may be, and should be allowed to be, risky and, therefore, may not work out. The ensuing work should be **subject to full and normal review.**"

*Software Engineering Research LABoratory*

## In general, would you participate again (as reviewer or authors)?

25 responses



- 🔵 Yes
- 🔴 No

100%

**Table 1** RR submissions and publications since inception at EMSE

| Venue | Stage 1 | | Stage 2 | |
| --- | --- | --- | --- | --- |
| | Submissions | IPAs | Submissions | Publications |
| MSR 2020 | 13 | 6 | 4 | 3 |
| MSR 2021 | 10 | 6 | 4 | 1 |
| MSR 2022 | 14 | 2 | 1 | 0 |
| ICSME 2020 | 7 | 4 | 3 | 2 |
| ICSME 2021 | n/a | 6 | 3 | 0 |
| ESEM 2021 | n/a | 4 | 1 | 0 |
| ESEM 2022 | 13 | 3 | 0 | 0 |

Note that some studies were affected by the COVID-19 pandemic. Data may be incomplete as tracking submissions can be challenging

# Open Issues and Questions

# Pros & Cons of RR

RRs provide early-stage feedback to authors and reduce researcher bias problems

**Table 2**  Benefits and disadvantages of registered reports in SE

| Benefits | Disadvantages |
|---|---|
| Shareable protocols for research replication. | More effort from researchers. |
| Focus is on research, not publication. | Limited acceptance by journals so far. |
| Improved rigour in reporting. | Rigour can mean different things to different people/communities (Storey et al. 2020). |
| Early peer review on research approach. | Not all research strategies are registerable. |

# Three faces of RR

**RR to prevent questionable research practices**

*Tell the world what you will do, then do it*

**RR as doctoral symposium**

*Early feedback before expensive data collection*

**RR as 1st round review**

*Pre-empt journal review with in-principle acceptance*

To what CS studies could it apply?

Most suited to post-positivist, confirmatory studies with clear hypotheses.

# Admin Challenges

CS has conference and journals - no one else does

Journals and conference **rarely share admin** interfaces (HotCRP vs Editorial Manager - and they are usually terrible)

Hard to manage reviewer discussions esp longitudinally

Currently, stick Phase 1 on Arxiv/OSF.io/Github

Have to **explicitly coach** reviewers (not yet mature, but true of other formats)

**Manually** track in progress RR on Google Sheets (low *vacation factor*)

# Admin Challenges

**Reviewer/editor burden** is increasingly a problem (overall, not just RR)

Accepting 5 IPAs at 3 conferences a year = 15 journal submissions in the next 12-18 months, with publication 24-36 months after that

+ who is asked to be conference track chair? What freedoms do they have?

**Minor shenanigans** - reviewer COI, authorship incentives

# Admin Challenges – J1C2?

Publication models run into journal profit models

First phase - Journal - then present at conference?

# RRs

**Enhance Reproducibility**

- Standardization of submitted protocols

**Are more likely to report Negative Results**

**Reviewers can help authors improve the protocol beforehand -> prevents flaws**

**Are a PLAN…. Not a PRISON**

- Flexibility is not lost … rather the possibility of airbrushing changes out of the picture

# Department of Reuse

Ultimately RR is about pre-specifying analysis. One way to do that is to reuse analysis protocols from other papers.

Done all the time in medicine; rarely in CS except in benchmarks.

Q: to what extent are artifacts such as protocols reused?



https://reuse-dept.org/

## Artifafct Creation, sharing and Reuse

**SE researchers share artifacts**

Not only publications …

Ideas, methods, datasets, tools

**Artifacts engage replication and reproducibility**

**Science produces more types of artifacts than just publications**

**Researchers use some but not not necessarily all artifacts from other work**

**HOW DO WE CAPTURE REUSE?**

43

# Badging – Artifact Evaluation Committees

The authors of accepted conference papers submit software packages that, in theory, let others re-execute that work.
These evaluation committees award "badges"

**Table 1. Badges such as the ones shown in this table are currently awarded at conferences.[2] This table is based on ACM's badge program, however, analogous badges are used at other conferences. Images used by permission of the Association for Computing Machinery.**

| Available | Functional | Reusable | Reproduced | Replicated |
|---|---|---|---|---|
| In a public repository with a long-term retention policy. A DOI needs to be provided. | Artifacts are documented, consistent, complete, exercisable, and include evidence of verification and validation. | Functional, significantly exceed minimal functionality. | Results of this paper have been reproduced by a different team using the original artifact. | Results of this paper have been replicated by a different team without the original artifact. |

# Badging – Artifact Evaluation Committees



Fig. 4. Artifact evaluation committee sizes 2011-2019. From Hermann et al. [8]

Is the artifact evaluation process is creating reused artifacts?

We queried ACM Portal for ICSE papers between 2011 to 2021, to find 2.4% of papers with an artifact badge.

Of these, 111 available, 74 reusable, 24 functional, NO replicated or reproduced artifacts.

approach to recording Research Reuse -> REUSE GRAPH

# Department of Reuse
- under development / data widely incomplete -

**Researchers**

Whose artifacts are reused (R+)

Jacques Klein
Yves Le Traon
Martin Monperrus
David Lo
Georgios Gousios

Who reuse artifacts (R)

Shangwen Wang
Xiaoguang Mao
Yang Liu
Ming Wen
Zhenchang Xing

**Artifacts**

Most reused (R+)

Towards Evaluating the Robustness of Neural Networks
Defects4J: a database of existing faults to enable controlled testing studies for Java programs
Deep Residual Learning for Image Recognition
Grounded theory in software engineering research
The GHTorrent dataset and tool suite

Most reusing (R)

Automated patch correctness assessment
Importance-driven deep learning system testing
Big code != big vocabulary
Debugging inputs
Fuzz testing based data augmentation to improve robustness of deep neural networks

**Statistics**

| | |
|---|---|
| Papers inspected | 128 |
| Reused papers (DOI) | 285 |
| Reused papers (arXiv) | |
| Reused repositories (GitHub) | |
| Reused websites | 77 |

Reuse types
Methodology (459), Tool (360), Dataset (193),
Metric (66), Stepping stone (57), Statistics (2
Sanity check (17), Replication (11), Other (6

**Legend**

● Inspected Paper
● Published Paper
● arXiv preprint
● GitHub Repository
● Grey Literature

⇨ Researchers read 170 SE papers selected from 6 major 2020 conferences

⇨ Teams were asked to record six types of reuse

⇨ Each edge connects papers to the prior work they are (re)using

This figure shows reuse from Bernal-Cárdenas et al.[7] Edges reflect tool, dataset, and methodology reuse. Red nodes indicate arXiv preprint; green represents a GitHub repository; blue denotes a published paper, and grey indicates other websites or grey literature locations. https://www.reuse-dept.org/doi/10.1145/3377811.3380328.

# ROSE festival (Rewarding Open Science Replication and Reproduction in SE)

**ICSE 2023**

## Call for Participation

Authors of papers with results that have been replicated or reproduced (*) by subsequent work (i.e. by **other** researchers) are invited to submit 1 one page ascii document to timm@ieee.org, title "ROSE'23 submission" that offers:

- a 4 line (or less) description of the original results
- a 4 line (or less) description of what was found by the other researchers
-  references to both the original paper and the subsequent work.

Accepted submissions will be offered a lightning talk slot at the ICSE'23 ROSE Festival.

DATES:
Submission: March 31, 2023
Notification: April 7: 2023
ROSE festival: dates TBD, some lunchtime in main ICSE conference

FOR MORE INFO:
timm@ieee.org

NOTES: (*)

### Important Dates — ⊕⊙ AoE (UTC-12h)

| | |
|---|---|
| **Fri 31 Mar 2023** Submission | new |
| **Fri 7 Apr 2023** Notification | new |

### ROSE

| | | |
|---|---|---|
| **Tim Menzies** North Carolina State University United States | Chair |
| **Neil Ernst** University of Victoria Canada | Chair |
| **Ben Hermann** TU Dortmund Germany | Chair |
| **Maria Teresa Baldassarre** Department of Computer Science, University of Bari Italy | Chair |

## Repeatability, Reproducibility, Replicability



| **Repeatability** | **Reproducibility** | **Replicability** |
|---|---|---|
| Original Team | Different Team | Different Team |
| Original Setup | Original Setup | Different Setup |

*Software Engineering Research LABoratory*

*The Rose Initiative (Recognizing and Rewarding Open Science in Software Engineering) is an international, multi-conference workshop that will continually report updates to the software engineering reuse graphs.*

# Credits & Special Thanks

**EDITORIAL**

## Registered reports in software engineering

Neil A. Ernst[1] · Maria Teresa Baldassarre[2]

**Abstract**
Registered reports are scientific publications which begin the publication process by first having the detailed research protocol, including key research questions, reviewed and approved by peers. Subsequent analysis and results are published with minimal additional review, even if there was no clear support for the underlying hypothesis, as long as the approved protocol is followed. Registered reports can prevent several questionable research practices and give early feedback on research designs. In software engineering research, registered reports were first introduced in the International Conference on Mining Software Repositories (MSR) in 2020. They are now established in three conferences and two pre-eminent journals, including this one (EMSE). We explain the motivation for registered reports, outline the way they have been implemented in software engineering, and outline some ongoing challenges for addressing high quality software engineering research.

**Keywords** Registered report · Research methods · Software engineering

## 1 Introduction

Registered reports are a model of scholarly publication which prioritize the importance of study design and significance rather than study outcomes. Focusing on whether the study was suitable to support the inferences of interest decouples publication from a focus on headline-worthy 'significant' results.

In software engineering (SE) research, empirical methods are now standard. The top conferences in the field emphasize "the extent to which the paper's contributions and/or

✉  Neil A. Ernst
    nernst@uvic.ca

    Maria Teresa Baldassarre
    mariateresa.baldassarre@uniba.it

[1]  Department of Computer Science, University of Victoria, Victoria, BC, Canada

[2]  Dipartimento di Informatica, Università degli studi di Bari, Bari, Italy

∅ Springer

# Credits & Special Thanks

# References

⇨ Chambers C. What's next for registered reports, Nature 573, 187-189 (2019) doi: https://doi.org/10.1038/d41586-019-02674-6

⇨ Popper, K. Conjectures and Refutations: The Growth of Scientific Knowledge. Routledge (1963)

⇨ Fabio Q. B. da Silva, Marcos Suassuna, A. César C. França, Alicia M. Grubb, Tatiana B. Gouveia, Cleviton V. F. Monteiro, and Igor Ebrahim dos Santos. 2012. Replication of empirical studies in software engineering research: a systematic mapping study. Empirical Software Engineering (Sept. 2012). https://doi.org/10.1007/s10664-012-9227-7

⇨ Ernst, N.A., Baldassarre, M.T. Registered reports in software engineering. Empir Software Eng 28, 55 (2023). https://doi.org/10.1007/s10664-022-10277-5

**Figure. Watch the authors discuss this work in the exclusive *Communications* video.**
https://cacm.acm.org/videos/reuse-of-research

mariateresa.baldassarre@uniba.it