Beyond the Badge: Reproducibility Engineering as a Lifetime Skill

Wolfgang Mauerer Technical University of Applied Science Regensburg Siemens AG, Corporate Research Germany wolfgang.mauerer@othr.de Stefan Klessinger Chair of Scalable Database Systems University of Passau Germany stefan.klessinger@uni-passau.de Stefanie Scherzinger Chair of Scalable Database Systems University of Passau Germany stefanie.scherzinger@uni-passau.de

ABSTRACT

Ascertaining reproducibility of scientific experiments is receiving increased attention across disciplines. We argue that the necessary skills are important *beyond* pure scientific utility, and that they should be taught as part of software engineering (SWE) education. They serve a dual purpose: Apart from acquiring the coveted badges assigned to reproducible research, reproducibility engineering is a *lifetime skill* for a professional industrial career in computer science.

SWE curricula seem an ideal fit for conveying such capabilities, yet they require some extensions, especially given that even at flagship conferences like ICSE, only slightly more than one-third of the technical papers (at the 2021 edition) receive recognition for artefact reusability. Knowledge and capabilities in setting up engineering environments that allow for reproducing artefacts and results over decades (a standard requirement in many traditional engineering disciplines), writing semi-literate commit messages that document crucial steps of a decision-making process and that are tightly coupled with code, or sustainably taming dynamic, quickly changing software dependencies, to name a few: They all contribute to solving the scientific reproducibility crisis, *and* enable software engineers to build sustainable, long-term maintainable, software-intensive, industrial systems. We propose to teach these skills at the undergraduate level, *on par* with traditional SWE topics.

CCS CONCEPTS

• Social and professional topics \rightarrow Software engineering education; • Software and its engineering \rightarrow Maintaining software; Software version control.

KEYWORDS

reproducibility engineering, teaching software engineering

1 INTRODUCTION

Since software engineering involves complex software stacks that non-trivially interact with hardware, sharing experimental setups is anything but trivial. Over the last decade, reproducibility of experimental results has become recognised as a prime aspect of computer science (CS) research. Several high-profile conferences now award badges when results can be independently verified.

Undoubtedly, reproducibility engineering (RepEng) has become a crucial skill that today's generation of PhD students has to master. In this position paper, we argue that these skills should already be taught (and practised) at the undergraduate level, and we therefore designed and conducted a course for computer science Bachelor students close to graduation. Even when students pursue an industry career, they will greatly benefit from recognising threats to reproducibility, how to tackle them, and how to build long-term reproducible code. In short, it is our conviction that students skilled in RepEng possess skills that proficient software engineers need to master (anyway).

Accordingly, we propose a multi-faceted syllabus¹ for teaching reproducibility engineering, and what we consider crucial skills. This includes best practices in computer science research and industry, such as packaging entire system software stacks for dissemination. For *long-term* reproducibility over decades (ideally, forever), we discuss why open source technologies (as massively employed in industry) are preferable to approaches crafted for research.

Structure. We recap essentials on building reproduction packages. We propose a high-level syllabus, covering social and technical best practices, as well as specific tools and technologies that are well-adopted in industry. We then discuss the literature material available for academic teaching, and conclude.

2 PRELIMINARIES

Terminology. Reproducibility is a cross-cutting theme and there are guidelines by the National Science Foundation (NSF) [17], the Association of Computing Machinery (ACM) [1], and the Institute of Electrical and Electronics Engineers (IEEE) [8]. Concepts like reproducibility and replicability receive different interpretation, depending on the community. Even large professional associations like the ACM had to revise their definitions of the terms because of prior confusion. Despite their obvious relevancy, the concepts are not yet reflected in the ACM Computing Classification System².

Throughout this article, we follow the ACM terminology [1] (version 1.1), and regard an experiment as *repeatable*, when the same team with the same experimental setup can confirm the results. An experiment is *reproducible*, if it is a different team, but the same setup, that confirms the results. Finally, an experiment is *replicable*, when a different team, with a different setup, confirms the result.

Reproduction Packages. Building a *reproduction package* goes beyond providing a document object identifier (DOI) to some repository hosting data, code, and setup instructions. Rather, a gold-standard reproduction package [7] bundles all research artefacts required to conduct the experiment (such as source code, libraries, or input

¹We have implemented the outlined ideas in an online course, taught in the winter term 2021/22, to undergraduate students at two universities. The lecture videos are available online on YouTube (link in PDF).

²Available online (link in PDF), last updated 2012.

data), and contains a dispatcher script that allows for executing and evaluating the experiment via a *single* command.



Figure 1 (adapted from [11]) shows a state-of-the-art setup. Based on system binaries, external and internal code in git repositories, and patch stacks with changes to existing components, a build recipe induces generation of a host-system independent Docker container as (static and immutable) build environment for measurement binaries (1). Additionally, a Docker container with pre-built binaries, devoid of any external dependencies, is created (2). The Docker container creates an experiment execution package (3) that can be deployed on cloud systems, or on local hardware, without any dependence on target-system provided artefacts. The experimental runs generate data, which are post-processed, evaluated and visualised by scripts in the experiment execution package.

3 A MULTI-LEVEL SYLLABUS

We argue that reproducibility engineering should find its way into undergraduate curricula, anchored in software engineering education. By targeting a clearly scoped audience (rather than STEM disciplines in general), we can address matters *to the point*, and provide actionable advice beyond the mechanical use of tools, or compliance to formal processes. In sketching out a syllabus, we propose a multi-level approach, and distinguish social and technical best practices. We further review specific tools and technologies.

3.1 Best practices: Social

A reproduction package should contain as many details as necessary, but must not overwhelm. Instead of minutiæ of how results were obtained, a reproduction package presents a concise and balanced view of the *outcome* of an effort. While any *structural decisions* are worth preserving, the temporal order of the *thought process* that led to intermediate results or to said decisions, is usually not.

Industry has established conventions [18] on documenting software changes (at the granularity of individual commits) to provide an understanding of the evolution of large software systems. These conventions can also be applied to documenting research progress. Such *trails of responsibility* (which persons authored changes together, who provided reviews, who participated in design decisions, etc.) are routinely created outside academia (contrariwise to the care taken in giving credit and attribution in published papers, this approach is not established in many areas of computer science). Figure 2 shows an example: It contains the technical change in form of a diff (bottom part), and metadata (unique hash, author and committer) as they are provided by version control systems like git. Apart from this information, as it is widely used in repository mining research [23], the commit also includes a summary of the change, and a rationale (brief for the sake of example) *why* the change is necessary, and *which* techniques are employed.

The commit can be seen as a form of communication with fellow humans instead of mere instructions for machines, following Knuth's seminal *literate programming* concept [10]. To create *readable histories*, we suggest to introduce the pragmatic customs developed in large, international and multi-disciplinary infrastructure projects (such as the Linux kernel) in software engineering courses.

commit: aa09c4f6a54152 ◀ Author: Jane Doe <jane@doe.com> ◀ Committer: John Doe <john@doe.com> ◀</john@doe.com></jane@doe.com>	- Unique ID of the commit - Author of change - Committer of change
Use salted hashes <	- Summary of changes
Function getHash() is used to hash user passwords. Since adding a salt value is considered a minimum standard these days, augment computing the hash with a salting function as devised by Ilsebill et al., Grassian Letters 27(3), 2022.	
Signed-off-by: Jane Doe <jane@doe.com> ← Credit for authorship Reviewed-by: Jean Doe <jean@doe.com> ← Credit for review Tested-by: Judy Doe <judy@doe.com> ← Credit for testing</judy@doe.com></jean@doe.com></jane@doe.com>	
diff -git a/sec/hash.c b/sec/hash.c ◀ @@ -1,7 +1,7 @@ doSomething();	——— Changed files
-hash = getHash(val); +hash = getSaltedHash(val, genSalt());	;



3.2 Best practices: Technical

We need to provide actionable guidance on how to implement the suggested procedures and approaches. This entails covering the necessary means *end-to-end*, from preparing all software components required to perform experiments, running analysis code and evaluations, and to creating insightful visualisations.

Building research artefacts depends on external sources, whose long-term availability is often not sufficiently considered. Particular care is taken to raise awareness for identifying potential issues when aiming at *reproducible builds*.

The granularity of packaging artefacts is an important discussion point: Should reproduction packages start directly with building the operating kernel from source, to establish absolutely identical conditions given identical hardware, or is it sufficient to package custom code that leverages any suitable execution platform? Likewise, should and can reproduction efforts re-compute all derived results, or start with data obtained from long-running calculations?

Another dimension concerns the variability of programming language, compiler and toolchain, and the distinction between build, execution, and evaluation platform. Each of the combinations that appear in practice have peculiarities worth discussing. Beyond the Badge: Reproducibility Engineering as a Lifetime Skill

Furthermore, we consider the technical ramifications of different types of reproducibility introduced in Section 2: Depending on what type of quantities are handled—physical quantities like time or energy consumption, numeric results from deterministic or stochastic processes, etc.—, different means ensure that it is possible to decide whether a reproduction attempt is successful.

Using closed-source, *proprietary components* creates hurdles for other researchers, and should be avoided in ideal open science. However, relying on proprietary components cannot be completely avoided, so we need to discuss how to best handle such scenarios.

Advanced numerical techniques that require accelerator hardware such as GPUs receive increasing attention in machine learning and artificial intelligence projects. The involved software stacks do not only contain binary-only, proprietary components whose licenses place obstacles on distribution, but also do not play well with virtualisation and containerisation approaches. We need to discuss how to handle these issues specific to *dealing with hardware*.

Finally, we need to address how to properly package artefacts and ensure their *long-term availability*. Besides using well-structured hierarchies and self-documenting package formats, we address dual strategies towards short- and long-term reproducibility: The latter aims at decades of reproducibility, at a higher cost to the reproducers, while the former accepts technologies and platforms that are not certified for *DOI-safety*, but allow for easier integration into standard development workflows. This balances advantages of long-term reproducibility with the ease of continuous development.

3.3 Tools and Technologies

We need to demonstrate tools that implement the previously discussed techniques. Primarily, we focus on Linux/Unix-based command-line tools, as these are also conveniently available on standard operating systems (Windows/Mac OS). This does not necessarily hold the other way (*e.g.*, Powershell), and for GUI approaches. A small subset of the tool functionality is sufficient for reproducibility engineering, and command-line based approaches are helpful locally and on servers. We suggest starting with means for easy, but effective, *low-threshold automation* based on efficient interaction with shells, pipelined processing of data, and *glue languages* such as python, R or Matlab.

Non-linear history rewriting provided by git allows for transforming chronological records into a readable, consistent research process documentation by splitting, merging, and re-ordering. Outside of software engineering, we have encountered little knowledge of such transformations, yet they are crucial to ascertain long-term human understandability.

Virtualisation and containers [2] play a major role in our strategy: For one, they avoid having to deal with different versions and compositions of compilers, libraries, and system software when building artefacts. Also, they allow for establishing a completely self-contained environment without external internet-based dependencies that remains operational even decades after the original sources have vanished (figuratively and literally, research is even possible when trapped on a remote island). Careful engineering of containers ensures they are suitable for reproduction tasks. The appropriate techniques and patterns should therefore be introduced. The reprotest tool collection is a recommended means to satisfy requirements for reproducible builds: By varying environmental parameters like user ID, folder names, or compiler settings, the tools detect issues that do not surface when a single researcher builds code on the always-same machine. Such setups lead, in our experience, to important insights on subtle sources of errors caused by implicit, yet common misconceptions. While such tools are routine for developers of distributions like Debian, and also key to long-term industrial maintainability of software, we find them not yet sufficiently integrated into SWE curricula.

We suggest to implement *preparing and documenting experiments* using knitr, which is not unsimilar to the paradigm of literate programming [5, 10]. It also allows for creating self-contained papers that realise end-to-end reproduction. *Electronic notebooks* like Jupyter are a recommended variant.

How to *describe and pin down the execution environment* is a further challenge. Typical hardware specifications in published research describe the experimental conditions along the lines of 'Linux version 5.1.92 on a Dull Powervortex 4711 with 24 GiB of RAM was used'. This is insufficient for reliable reproduction—non-standard kernel extensions that may vary widely depending on the distribution, specific settings for tuning parameters that exist in a wide variety on every system, and many other factors that may easily be dismissed as irrelevant can impact measurements by orders of magnitude. We recommend discussing means of faithfully recording the execution conditions of computational experiments.

Students should acquire hands-on experience in *reproducing experimental outcomes*. Retracing the work of others increases awareness for (and appreciation of) high-quality reproduction packages.

3.4 Special Cases

While software engineering can often be decoupled from details of the target environment (CPU architecture, OS version, ...), specialpurpose hardware introduces additional reproducibility requirements. We find that general-purpose graphical processing units (GPGPUs) necessitate software stacks that exceed standard compilers considerably in size, and introduce (a) strong interdependencies between software component versions and (b) dependencies on specific features that might only find intermittent support in hardware. Both stress the need to teach implementing less performant, but generic alternatives, and how to store intermediate results obtained from HW accelerators for further processing. Similar considerations hold for tensor processing units (TPUs) and other AI accelerators. Quantum computing, starting to receive interest from the software engineering community [19, 22], additionally needs to deal with globally unique hardware semi-prototypes [14].

4 TEACHING MATERIAL

Textbooks. While a number of books on the subject itself have been published, they are either (a) edited collections of articles written by different authors and lack a central *leitmotif*, (b) focus on high-level aspects of reproducibility, or (c) consider very narrow domains. For instance, recent books discuss reproducibility in preclinical animal studies [16], biomedical sciences [15, 21], or pattern recognition [9], with limited applicability outside these fields. The book by Stodden, Leisch and Peng [20] comes, despite being a collection of articles, close to what we need in academic teaching: it seeks to augment general advice on reproducibility engineering with concrete technical details and examples. However, almost all of the recommended tools—with the notable exception of knitr, which we also include in our recommendations—stem from scientific research. At the time of this writing, they are no longer maintained (VisTrail [4]), fail to build (Sumatra), or are no longer available, apart from historical archives like the wayback machine (CDE, SOLE). Given that the book was published in 2014, this underlines our strategy of relying on industrial, long-term maintained tools, as academic tools tend to break once project funding ceases [6].

Online courses. Several MOOC platforms offer courses on reproducibility engineering (*e.g.*, Coursera, EdX, Inria). They cover topics related to software engineering (such as literate programming via notebooks), but originate from outside the SWE community, (*e.g.*, computational biology or biomathematics). Our own online course (cf. Sec. 1) assumes the computer science perspective.

5 SUMMARY, EXPERIENCE AND OUTLOOK

Reproducibility engineering prepares students towards industry careers, where sustainable long-term maintenance is important. It should also become an entry-level requirement for PhD candidates.

In teaching and evaluating the class, we have observed that these goals were satisfied. Difficulties arise when technical details subtly impact reproducibility (*e.g.*, different CPU architecture between VM and host, or host CPU details). A solution was to often add additional packages and layers instead of identifying root causes of non-reproducibility. Consequently, we find that details may matter to a larger extent than in other aspects of software engineering.

Finally, we believe the effort contributes towards solving the reproducibility crisis. Computer science and software engineering seem, in a pivotal function, predestined for this purpose.

Acknowledgements. The joint effort was partly funded by the Lehrinnovationspool 2.0/2021-2022 at the University of Passau. The authors were partly funded by Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) grant #385808805, and BMBF grant number 13N15645. We thank Pia Eichinger for assisting with the video takes.

AUTHORS' PROFILES



WOLFGANG MAUERER (rhs) is a professor at Technical University of Applied Sciences Regensburg, and a senior research scientist at Siemens AG, Corporate Research. STEFANIE SCHERZINGER (lhs) is a professor at University of Passau, where

she chairs the *Scalable Database Systems* group. Together, they have taught several courses and tutorials [12, 13], including the inverted classroom on reproducibility engineering featured here. They have also carried out reproduction studies of published research [3].



STEFAN KLESSINGER, M. Sc., is a computer science researcher at University of Passau, and has tutored the joint RepEng course.

REFERENCES

- ACM 2020. Review and Badging Artifact. https://www.ieee.org/publications/ research-reproducibility.html [Online].
- [2] Carl Boettiger. 2015. An Introduction to Docker for Reproducible Research. SIGOPS Oper. Syst. Rev. 49, 1 (Jan. 2015), 71–79. https://doi.org/10.1145/2723872. 2723882
- [3] Dimitri Braininger, Wolfgang Mauerer, and Stefanie Scherzinger. 2020. Replicability and Reproducibility of a Schema Evolution Study in Embedded Databases. In Proc. EmpER 2020. 210–219. https://doi.org/10.1007/978-3-030-65847-2_19
- [4] Steven P. Callahan, Juliana Freire, Émanuele Santos, Carlos E. Scheidegger, Cláudio T. Silva, and Huy T. Vo. 2006. VisTrails: Visualization Meets Data Management. In Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data (Chicago, IL, USA) (SIGMOD '06). Association for Computing Machinery, New York, NY, USA, 745–747. https://doi.org/10.1145/1142473.1142574
- [5] Jon F. Claerbout and Martin Karrenbach. 2005. Electronic documents give reproducible research a new meaning. 601–604. https://doi.org/10.1190/1.1822162 arXiv:https://library.seg.org/doi/pdf/10.1190/1.1822162
- [6] Sergey Fomel. 2015. Reproducible Research as a Community Effort: Lessons from the Madagascar Project. Computing in Science and Engg. 17, 1 (Jan. 2015), 20-26. https://doi.org/10.1109/MCSE.2014.94
- [7] Benjamin J. Heil, Michael M. Hoffman, Florian Markowetz, Su-In Lee, Casey S. Greene, and Stephanie C. Hicks. 2021. Reproducibility standards for machine learning in the life sciences. *Nature Methods* 18, 10 (Aug. 2021), 1132–1135. https://doi.org/10.1038/s41592-021-01256-7
- [8] IEEE 2016. Report on the First IEEE Workshop on The Future of Research Curation and Research Reproducibility. https://www.ieee.org/publications/researchreproducibility.html [Online].
- [9] Bertrand Kerautret, Adrien Krähenbühl, Pascal Monasse, Miguel Colom, Daniel Lopresti, and Hugues Talbot. 2021. Reproducible Research in Pattern Recognition. Springer. https://doi.org/10.1007/978-3-030-76423-4
- [10] Donald E. Knuth. 1984. Literate Programming. Comput. J. 27, 2 (1984), 97–111. http://dblp.uni-trier.de/db/journals/cj/cj27.html#Knuth84
- [11] Wolfgang Mauerer, Ralf Ramsauer, Edson Ramiro Lucas Filho, and Stefanie Scherzinger. 2021. Silentium! Run-Analyse-Eradicate the Noise out of the DB/OS Stack. In Proc. Fachtagung für Datenbanksysteme für Business, Technologie und Web (BTW) 2021. https://doi.org/10.18420/btw2021-21
- [12] Wolfgang Mauerer and Stefanie Scherzinger. 2020. Educating Future Software Architects in the Art and Science of Analysing Software Data. In Proc. "Software Engineering im Unterricht der Hochschulen" 2020 (CEUR Workshop Proceedings, Vol. 2531). CEUR-WS.org, 56–60.
- [13] Wolfgang Mauerer and Stefanie Scherzinger. 2021. Nullius in Verba: Reproducibility for Database Systems Research, Revisited. In 37th IEEE International Conference on Data Engineering, ICDE 2021. IEEE, 2377–2380. https: //doi.org/10.1109/ICDE51399.2021.00270
- [14] Wolfgang Mauerer and Stefanie Scherzinger. 2022. 1-2-3 Reproducibility for Quantum Software Experiments. In In Proc. 1st International Workshop on Quantum Software Analysis, Evolution and Reengineering (Q-SANER@SANER 2022).
- [15] Erwin B Montgomery Jr. 2019. Reproducibility in Biomedical Research: Epistemological and Statistical Problems. Academic Press.
- [16] José M. Sánchez Morgado and Aurora Brønstad (Eds.). 2021. Experimental Design and Reproducibility in Preclinical Animal Studies. Springer.
- [17] National Academies of Sciences, Engineering, and Medicine and others. 2019. Reproducibility and replicability in science. National Academies Press.
- [18] Ralf Ramsauer, Daniel Lohmann, and Wolfgang Mauerer. 2019. The list is the process: reliable pre-integration tracking of commits on mailing lists. In Proc. ICSE 2019. IEEE / ACM, 807–818. https://doi.org/10.1109/ICSE.2019.00088
- [19] Manuel Schönberger, Maja Franz, Stefanie Scherzinger, and Wolfgang Mauerer. 2022. Peel | Pile? Cross-Framework Portability of Quantum Software. QSA@ICSA 2022.
- [20] Victoria Stodden, Friedrich Leisch, and Roger D. Peng (Eds.). 2014. Implementing Reproducible Research. CRC Press.
- [21] Michael Williams, Michael Curtis, and Kevin Mullane. 2017. Research in the biomedical sciences: Transparent and reproducible. Academic Press.
- [22] Jianjun Zhao. 2020. Quantum software engineering: Landscapes and horizons. (2020). arXiv:2007.07047
- [23] Thomas Zimmermann, Andreas Zeller, Peter Weissgerber, and Stephan Diehl. 2005. Mining version histories to guide software changes. Software Engineering, IEEE Transactions on 31, 6 (2005), 429–445.