

Artifact Evaluations as Authors and Reviewers: Lessons, Questions, and Frustrations

2024 Community Workshop on Practical Reproducibility in HPC

Quentin Guilloteau¹, Millian Poquet², Jonas H. Müller Korndörfer¹, Florina M. Ciorba¹

2024-11-18

¹University of Basel, DMI, HPC, Basel, Switzerland

²University of Toulouse, IRIT, Toulouse, France



University
of Basel



DAPHNE



UNIVERSITÉ
TOULOUSE III
PAUL SABATIER



Who are we?

Research Topics

- HPC
- Scheduling (OMP, MPI, RJMS, etc.)
- **Reproducible Research!**

Recent Activities and AE Experiences

- Artifact Reviewers (SC24, EuroSys'25)
- Artifact Authors (Euro-Par24, TPDS'22)
- Attendees and Organizers of Reproducibility Hackathons (<https://www.reprohack.org/>)
- **Study of ADs in HPC (ACM REP24)**



Longevity of Artifacts in Leading Parallel and Distributed Systems Conferences: a Review of the State of the Practice in 2023

Quentin Guilloteau
Florina M. Cioba
Quentin.Guilloteau@univ-chu.fr
Florina.Cioba@univ-chu.fr
University of Basel
Basel, Switzerland

Millian Poquet
Millian.Poquet@univ-tlse.fr
Univ. Toulouse, CNRS, IRIT
Toulouse, France

Dorian Goepff
Dorian.Goepff@univ-grenoble.fr
Dorian.Goepff@univ-grenoble.fr
Univ. Grenoble Alpes, Inria, CNRS, IIG
Grenoble, France

ABSTRACT

Reproducibility in the computer science. Many scientific communities have been struck by the reproducibility crisis, and computer science is no exception. Its answer has been to require artifact evaluations along with accepted articles and award badges to reward authors for their efforts to support reproducibility. Authors voluntarily submit artifacts associated with a submission to reviewers who decide their "reproducibility" properties. We report the results of "reproducibility" considered by such badges (visited and reused) important aspects of the reproducibility crisis. In this article, we survey almost 1000 articles from five leading conferences on parallel and distributed systems held in 2023 (ACM SIGPLAN, EuroSys, OSDI, PPoPP, and SC). For each article, we gather information about its artifacts (how it was shared, under which experimental setup, and how the software environment was generated and shared), as well as the reproducibility badges awarded. By reviewing the state-of-practice does not address reproducibility in terms of artifact longevity and we expose eight observations that support this finding. To address the longevity of artifacts, we propose a new badge based on source code, experimental setup, and software environment. These criteria will allow reviewing artifacts exposed to software test of time. This work aims to shed light on the issue of long reproducibility in parallel and distributed systems and to discuss in the community towards addressing the issue.

CCS CONCEPTS

General and reference → Empirical studies.

KEYWORDS

Reproducibility, Artifact Evaluation, Badges, Longevity

ACM Reference Format:

Quentin Guilloteau, Florina M. Cioba, Millian Poquet, Dorian Goepff, and Olivier Richard. 2023. Longevity of Artifacts in Leading Parallel and Distributed Systems Conferences: a Review of the State of the Practice in Its ACM Conference on Reproducibility and Repeatability (ACM REP).



This work is licensed under a Creative Commons Attribution 4.0 International License.

ACM REP '23, June 18–20, 2023, Rennes, France.
© 2023 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-1023-3/23/06.
<https://doi.org/10.1145/3541233.3541233>

June 18–20, 2023, Rennes, France, ACM, New York, NY, USA, 13 pages.
<https://doi.org/10.1145/3541233.3541233>

1 INTRODUCTION

The scientific community as a whole is traversing a reproducibility crisis for the last decade. Computer science is not an exception to this crisis [4, 65]. The reproducibility of research is essential to build solid knowledge and increase reliability and confidence in the results, while limiting the methodology and analysis bias. In 2015, Collberg et al. [13] studied the reproducibility of 402 experimental articles published in system conferences and journals of 2011 and 2012. Each of the articles studied linked the source code used to perform their experiments. Of the 402 articles, 44% were not reproducible. The main causes were: (i) the source code was not available, (ii) the code did not compile or run, (iii) the experiments required specific hardware.

To reward authors of reproducible articles, several publishers, such as ACM or Springer, set up a peer review-based artifact evaluation for each submission. This peer review process of the experimental artifact can reward one or several badges to the authors based on the level of reproducibility of their artifacts.

The term reproducibility is often used as a broad name and



Who are we?

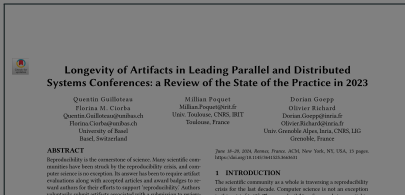
Research Topics

- HPC
- Scheduling (OMP, MPI, RJMS, etc.)
- **Reproducible Research!**

Rec

- Artifact Reviewers (SC21, EuroSys 23)
- Artifact Authors (Euro-Par24, TPDS'22)
- Attendees and Organizers of Reproducibility Hackathons (<https://www.reprohack.org/>)
- **Study of ADs in HPC (ACM REP24)**

This presentation: Feedback from all our experiences



The AE Process in a Nutshell

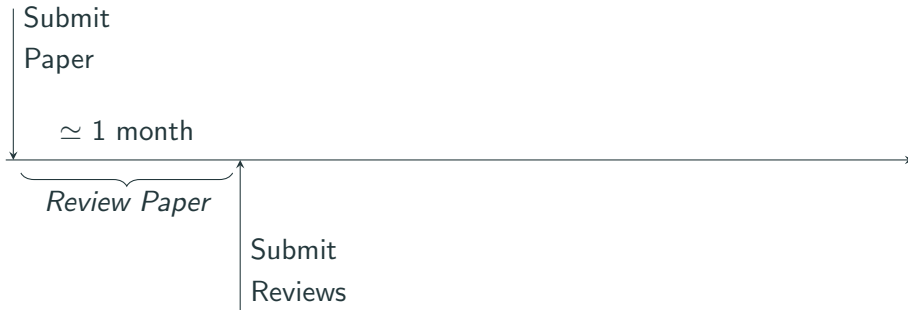
Authors



Timeline not to scale!

The AE Process in a Nutshell

Authors

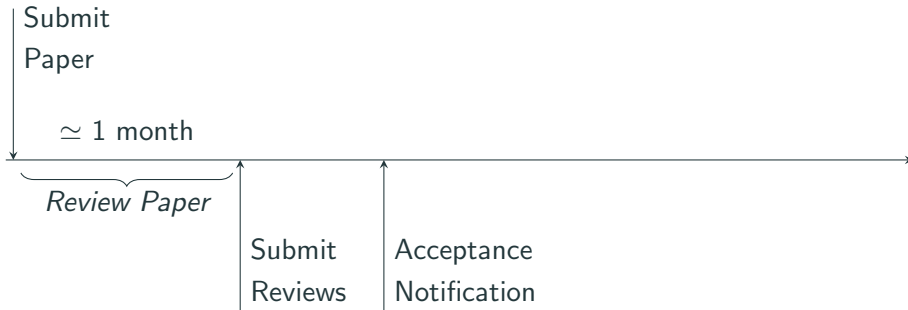


Reviewers

Timeline not to scale!

The AE Process in a Nutshell

Authors



Reviewers **Chairs**

Timeline not to scale!

The AE Process in a Nutshell

Authors

Authors

Submit
Paper

Submit
Artifact

≈ 1 month

Review Paper

Submit
Reviews

Acceptance
Notification

Reviewers

Chairs

Timeline not to scale!

Appendix: Artifact Description/Artifact Evaluation

Artifact Description (AD)

I. OVERVIEW OF CONTRIBUTIONS AND ARTIFACTS

A. Paper's Main Contributions

Provide a list of all main contributions of the paper.

- C₁ This is the 1st contribution.
- C₂ This is the 2nd contribution.
- C₃ This is the 3rd contribution.

B. Computational Artifacts

List the computational artifacts related to this paper along with their respective DOIs. Note that all computational artifacts may be archived under a single DOI.

- A₁ <https://doi.org/YYYYY/nomodo.0XXXXX>
- A₂ <https://doi.org/ZZ/YYYY/nomodo.1XXXXX>
- A₃ <https://doi.org/ZZ/YYYY/nomodo.2XXXXX>

Provide a table with the relevant computational artifacts, highlight their relation to the contributions (from above) and point to the elements in the paper that are reproducible by each artifact, e.g., which figures or tables were generated with the artifact.

Artifact ID	Contributions Supported	Related Paper Elements
A ₁ , C ₁	Tables 1-2	Figure 3
A ₂ , C ₂	Tables 2-3	Figures 1-2
-	-	-

II. ARTIFACT IDENTIFICATION

Provide the following six subsections for each computational artifact A_i.

A. Computational Artifact A_i

Relation to Contributions

Briefly explain the relationship between the artifact and contributions.

Expected Results

Provide a higher level description of what outcome to expect from the corresponding experiments. Provide an explanation of how the results substantiate the main contributions.

Algorithm A should be faster than Algorithms C and B in all GPU scenarios.

Expected Reproduction Time (in Minutes)

Estimate the time required to reproduce the artifact, providing separate estimates for the individual steps: Artifact Setup, Artifact Execution, and Artifact Analysis.

The expected computational time of this artifact on GPU X is 20 min.

Artifact Setup (incl. Inputs)

Hardware: Specify the hardware requirements and dependencies (e.g., a specific internet or GPU type is required).

Software: Introduce all required software packages, including the computational artifact. For each software package, specify the version and provide the URL.

Parameters / Inputs: Describe the datasets required by the artifact. Indicate whether the datasets can be generated, including instructions, or if they are available for download, providing the corresponding URL. **Installation and Deployment:** Detail the requirements for compiling, deploying, and executing the experiments, including necessary compilers and their versions.

Artifact Execution

Provide an abstract description of the experiment workflow of the artifact. It is important to identify the main tasks (processors) and how they depend on each other.

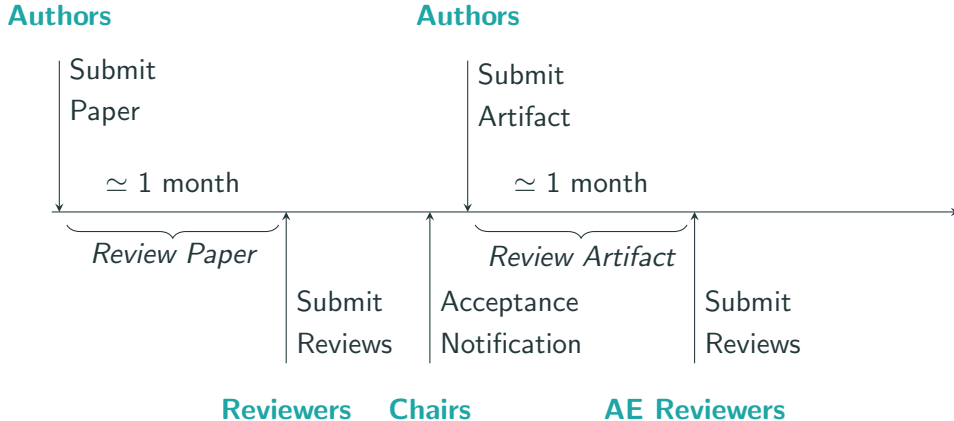
A workflow may consist of three tasks: T₁, T₂, and T₃. The task T₁ may generate a specific dataset. This dataset is then used as input by a computational task T₂, and the output of T₂ is processed by another task T₃, which produces the final results (e.g., plots, tables, etc.). Show the individual tasks T_i and provide their dependencies, e.g., T₁ → T₂ → T₃.

Provide details on the experimental parameters. How and why were parameters set to a specific value (if relevant for the reproduction of an artifact), e.g., size of dataset, number of data points, input sizes, etc. Additionally, include details on statistical parameters, like the number of repetitions.

Artifact Analysis (incl. Outputs)

B. Computational Artifact A_i
Provide the same type of information as done for Computational Artifact A₁.

The AE Process in a Nutshell



Timeline not to scale!

The AE Process in a Nutshell

Authors

Submit
Paper

\simeq 1 month

Review Paper

Submit
Reviews

Authors

Submit
Artifact

\simeq 1 month

Review Artifact

Acceptance
Notification



Submit
Reviews

Award
Badges

Reviewers

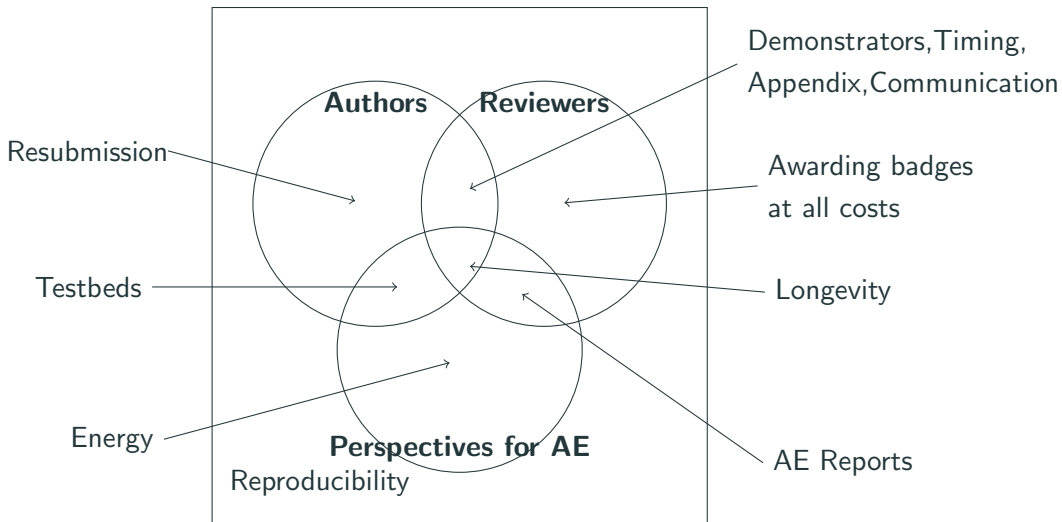
Chairs

AE Reviewers

AE Chairs

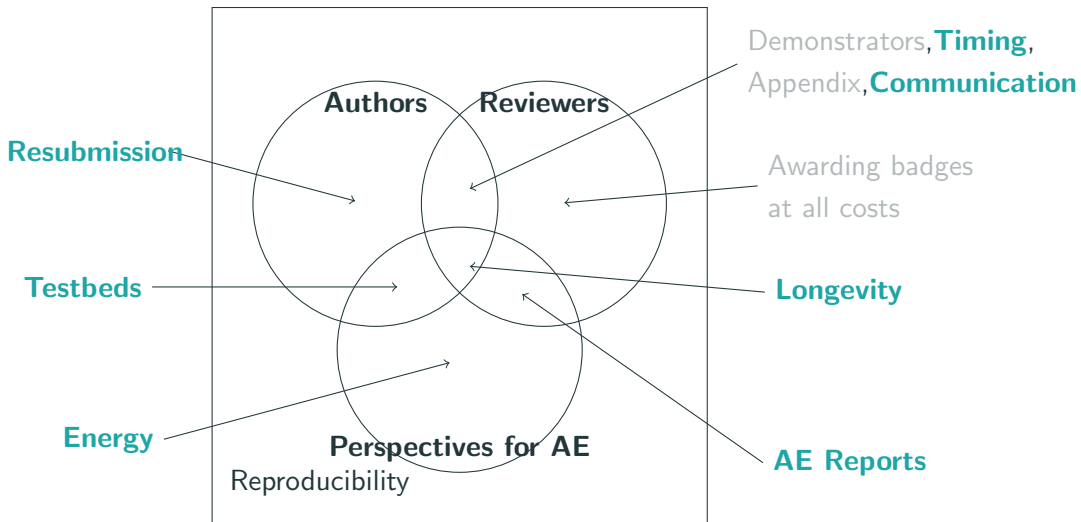
Timeline not to scale!

Topics Addressed in this Presentation



(If you were/are an AE Chair: share your perspective!)

Topics Addressed in this Presentation



(If you were/are an AE Chair: share your perspective!)

Timing of AE

- Artifacts **quickly** created by the authors close to the deadline
- AE Reviewers have a short time to evaluate in a **busy** schedule
- **Available** physical resources might be hard to find during evaluation process

Timing of AE

- Artifacts **quickly** created by the authors close to the deadline
- AE Reviewers have a short time to evaluate in a **busy** schedule
- **Available** physical resources might be hard to find during evaluation process

↪ Is the AE process **too short / rushed?**

Timing of AE

↔ Is the AE process **too short / rushed?**

Communication Between Authors and AE Reviewers

Timing of AE

↪ Is the AE process **too short / rushed**?

Communication Between Authors and AE Reviewers

↪ What is the role of AE Reviewers: **debug or evaluate** artifacts?

Timing of AE

↔ Is the AE process **too short / rushed**?

Communication Between Authors and AE Reviewers

↔ What is the role of AE Reviewers: **debug or evaluate** artifacts?

Resubmission Not Possible

Timing of AE

↪ Is the AE process **too short / rushed**?

Communication Between Authors and AE Reviewers

↪ What is the role of AE Reviewers: **debug or evaluate** artifacts?

Resubmission Not Possible

↪ If no badge, then **“artifact rejection”** \rightsquigarrow no resubmission

Timing, Communication, and Resubmission

Timing of AE

↪ Is the AE process **too short / rushed**?

Communication Between Authors and AE Reviewers

Does the AE process need to be **redesigned**?
Only for journals? Required **before** submission?

Resubmission Not Possible

↪ If no badge, then “**artifact rejection**” \rightsquigarrow no resubmission

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)
- Recent Findings:
 - Poor control of artifact sources
(dead links, commit not fixed/specified)
↪ **Zenodo, Software-Heritage**

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)
- Recent Findings:
 - Poor control of artifact sources
(dead links, commit not fixed/specified)
↪ **Zenodo, Software-Heritage**
 - Poor control of software environments
(package list, sometimes versions, apt, pip, userspace only)
↪ **Nix(OS)/Guix(System)**

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)
- Recent Findings:
 - Poor control of artifact sources
(dead links, commit not fixed/specified)
↪ **Zenodo, Software-Heritage**
 - Poor control of software environments
(package list, sometimes versions, apt, pip, userspace only)
↪ **Nix(OS)/Guix(System)**
 - Hardware not easy to access
↪ **Experimental testbeds**

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)
- Recent Findings:
 - Poor control of artifact sources
(dead links, commit not fixed/specified)
⇒ **Zenodo, Software-Heritage**
 - Poor control of software environments
(package list, sometimes versions, apt, pip, userspace only)
⇒ **Nix(OS)/Guix(System)**
 - Hardware not easy to access
⇒ **Experimental testbeds**



New Badge?

Longevity

- Science: **Iterative Self-Correcting Process**
(“standing on the shoulders of giants” – Isaac Newton)
- Recent Findings:
 - Poor control of artifact sources
(dead links, commit not fixed/specified)
↪ **Zenodo, Software-Heritage**
 - Poor control of software environments
(package list, sometimes versions, apt, pip, userspace only)
↪ **Nix(OS)/Guix(System)**
 - Hardware not easy to access
↪ **Experimental testbeds**



New Badge?

↪ **Who is the target** of reproducible research?

Longevity

↪ **Who is the target** of reproducible research?

AE Reports

Longevity

↪ **Who is the target** of reproducible research?

AE Reports

- Badges \simeq pass/fail
- Review reports might contain complementary information to AD and badges
↪ **no access** for future researchers
- AE Reports allow for **more details** about the reviewers' reproduction attempts

Longevity

↪ **Who is the target** of reproducible research?

AE Reports

- Badges \simeq pass/fail
- Review reports might contain complementary information to AD and badges
↪ **no access** for future researchers
- AE Reports allow for **more details** about the reviewers' reproduction attempts

↪ Do badges carry **enough nuance/information**?

Longevity

↪ **Who is the target** of reproducible research?

AE Reports

Who should benefit from AE?
Authors and/or future researchers?

↪ **no access** for future researchers

- AE Reports allow for **more details** about the reviewers' reproduction attempts

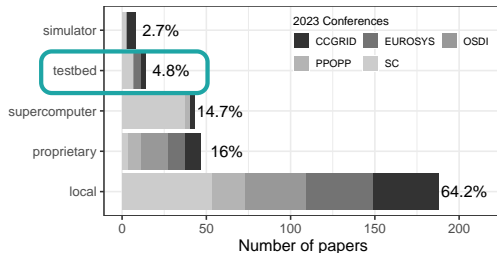
↪ Do badges carry **enough nuance/information?**

Testbeds (Chameleon, CloudLab, Grid'5000, etc.) are **amazing** to access hardware and reduce unknown/variability of deployment.

Testbeds (Chameleon, CloudLab, Grid'5000, etc.) are **amazing** to access hardware and reduce unknown/variability of deployment.

But...

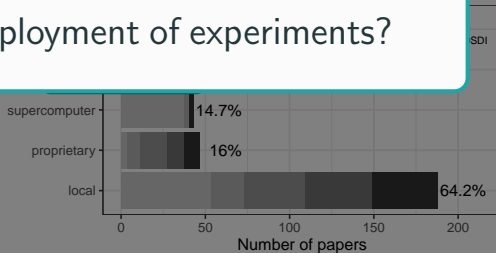
- Waiting time $>$ AE Reviewing time
- Not widely used in practice by authors
- Strong dependency on the testbed?



Testbeds (Chameleon, CloudLab, Grid'5000, etc.) are **amazing** to access hardware and reduce unknown/variability of deployment.

How to **“future-proof”** the deployment of experiments?

- waiting time > AC reviewing time
- Not widely used in practice by authors
- Strong dependency on the testbed?



HPC experiments consume **time and energy**

- (and increasingly more of each with AI in HPC conferences)
- Should we store **everything**?
 - why storing the result instead of the recipe?
 - \leftrightarrow need reproducible/deterministic ways to produce research objects
 - Nix/Guix? but might recompile *a lot* \rightsquigarrow Sustainable?
- Minimal viable example/experiment: but must be **representative** of the study
- How to create a valuable minimal viable example/experiment?
- How to reward a partially evaluated artifact?
- Is it worth to ask **several reviewers** to try to reproduce all or part of the study?

HPC experiments consume **time and energy**

- (and increasingly more of each with AI in HPC conferences)
- Should we store **everything**?

When do the **time and energy costs outweigh the value** of what is reproduced?

- Minimal viable example/experiment: but must be **representative** of the study
- How to create a valuable minimal viable example/experiment?
- How to reward a partially evaluated artifact?
- Is it worth to ask **several reviewers** to try to reproduce all or part of the study?

Is the Artifact Evaluation process (creation + evaluation) **rushed**?

Is **Artifact Evaluation** the path to **Reproducibility** in HPC?

Who is the target of reproducible research?

What is the future of AE in HPC in an **energy-constrained** world?