

Longitudinal Study of Software Environments Produced by Dockerfiles from Research Artifacts: Initial Design

ACM REP'25, Vancouver, BC, Canada

Quentin GUILLOTEAU, Antoine WAEHREN, Florina M. CIORBA

2025-07-30

University of Basel, Switzerland

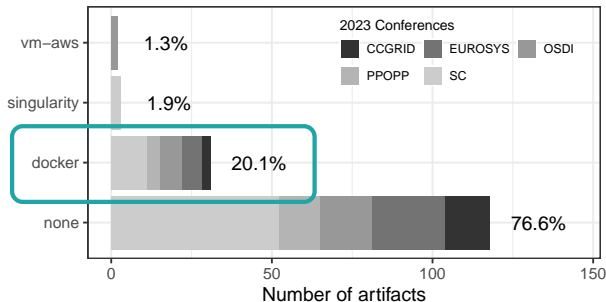


University
of Basel

Reproducibility, Artifacts, and Longevity

Longevity (ACM REP'24)

- Who are artifacts for? → **Future researchers**
- Math proof will not **disappear** or **change**
- Conferences recommend using containers



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

Quentin Guillemin, Florian M. Ciorba, Quentin Guillemin@uniba.ch, Florian.Ciorba@uniba.ch, University of Basel, Basel, Switzerland

Millian Piquet, Millian.Piquet@uniba.ch, Univ. Toulouse, CNRS, IRIT, Toulouse, France

Dorian Goepp, Olivier Richard, Dorian.Goepp@uniba.ch, Olivier.Richard@uniba.ch, Univ. Grenoble Alpes, CNRS, IIG, Grenoble, France

ABSTRACT

Reproducibility is the cornerstone of 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 argue that the notion of "reproducibility" considered by such badges is limited and raises important aspects of the reproducibility crisis. In this article, we survey almost 300 articles from five leading conferences on parallel and distributed systems held in 2023 (CCGrid, 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 methods and tools used to create and share artifacts in a holistic, in-depth, and article-context-specific manner, we found that the state-of-the-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 rewarding artifacts expected to withstand test of time. This work aims to shed light on the issue of long reproducibility in parallel and distributed systems and to its discussion in the community towards addressing the issue.

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, 45]. 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 2005, Collopy et al. [15] studied the reproducibility of 402 experimental articles published in system conferences and journals of 2001 and 2002. Each of the articles studied linked the source code used to perform their experiments. Of the 402 articles, 46% 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 award one or several badges to the authors based on the level of reproducibility of their artifacts.

CCS CONCEPTS

General and reference → Empirical studies

KEYWORDS

Reproducibility, Artifact Evaluation, Badges, Longevity

ACM Reference Format:

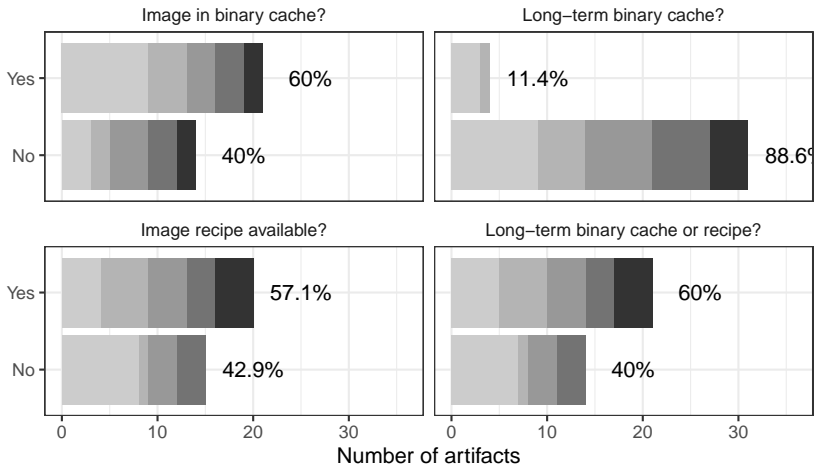
Quentin Guillemin, Florian M. Ciorba, Millian Piquet, Dorian Goepp, and Olivier Richard. 2023. Longevity of Artifacts in Leading Parallel and Distributed Systems Conferences: a Review of the State of the Practice in 2023. In *ACM Conference on Reproducibility and Replicability (ACM REP)*.

BY

Permission is granted under a Creative Commons Attribution International 4.0 License.

ACM REP '24, June 10–20, 2024, Annecy, France.
© 2024 Copyright held by the owner/authors.
ACM ISBN 978-1-4503-4234-6.
<https://doi.org/10.1145/3641233.3643031>

Docker, Longevity, and Sustainability



2023 Conferences



CCGRID



EUROSYS



OSDI

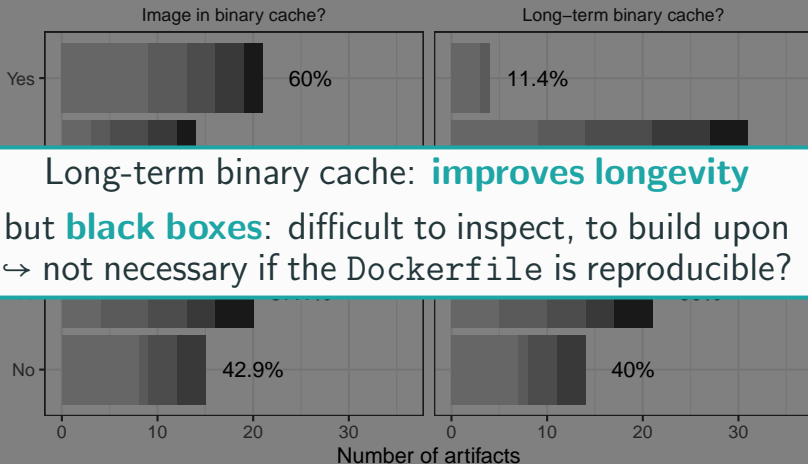


PPOPP



SC

Docker, Longevity, and Sustainability



2023 Conferences

CCGRID

EUROSYS

OSDI

PPOPP

SC

Is Docker really suitable for **longevous and reproducible** research?

(Should Reproducibility Chairs stop recommending Docker as a suitable solution to authors?)

Is Docker really suitable for **longevous and reproducible** research?

(Should Reproducibility Chairs stop recommending Docker as a suitable solution to authors?)

↪ How do the software environments produced by Dockerfiles
evolve through time?

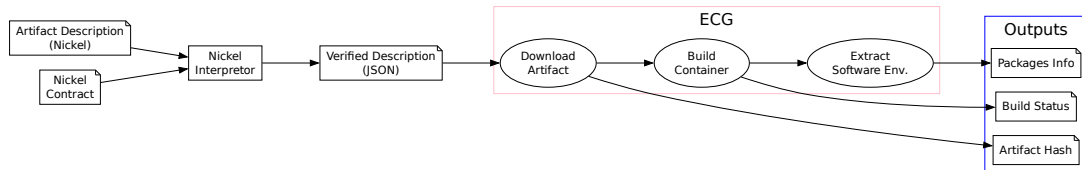
Is Docker really suitable for **longevous and reproducible** research?

(Should Reproducibility Chairs stop recommending Docker as a suitable solution to authors?)

Study: Take Dockerfiles from Research Artifacts, **build** them periodically, and **capture** the resulting software environment.

→ How do the software environments produced by Dockerfiles **evolve through time**?

Workflow, Data collected, and Frequency



What are we capturing?

- **Artifact Hash**: Did the content of the artifact change? (what's behind the link)
- **Build Status**: Did the container build successfully? What were the errors?
- **Packages Info**: What are the versions of the packages in the SW environment?
 - ↪ **Package Managers** (apt, dpkg, pip, conda), **Manual Installs.** (git, curl/wget)

When are we capturing? And for how long?

At the start of **each month**, for a **full year** (13 captures)

Scope of this Preliminary Study

5 artifacts from Euro-Par 2025 \leadsto all the artifacts using a Dockerfile

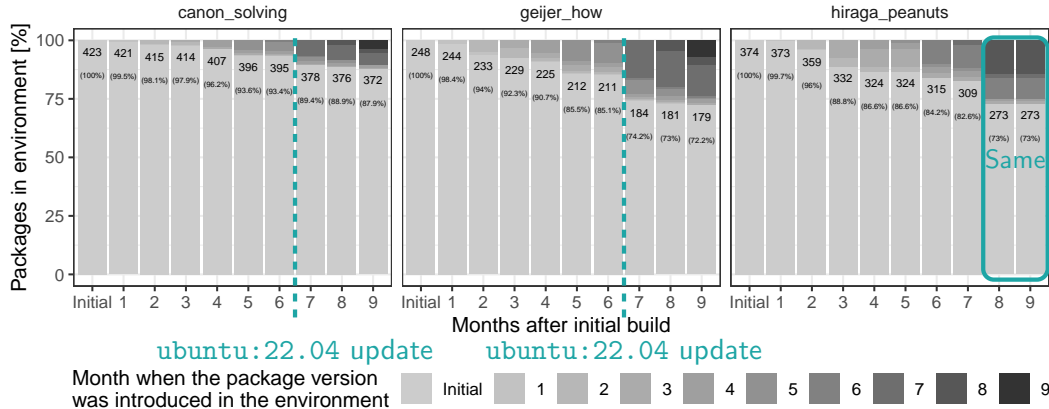
- Conference in our field (HPC) \leadsto **we are familiar with SW stacks**
- Artifacts published when we finished to develop our workflow \leadsto **“fresh” artifacts**

Artifact	Docker Base Image used		Calling apt update?
	Name	Version	
canon_solving	ubuntu	22.04	Yes
geijer_how	ubuntu	22.04	Yes
hiraga_peanuts	devcontainers/cpp	1-debian-12	Yes
munoz_fault	ubuntu	22.04	Yes
wolff_fast	ubuntu	22.04	Yes

Table 1: Information about the Dockerfiles from the study.

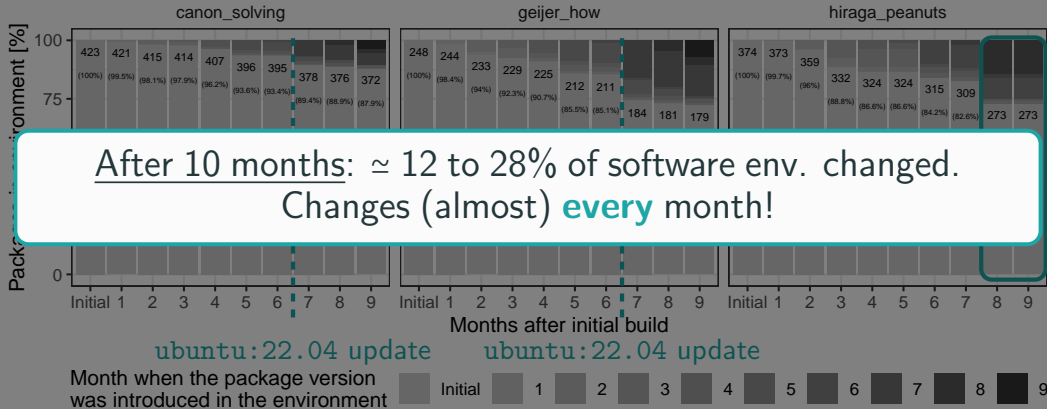
Preliminary Results – Per Artifact

Evolution of the packages versions over time for selected Dockerfile studied



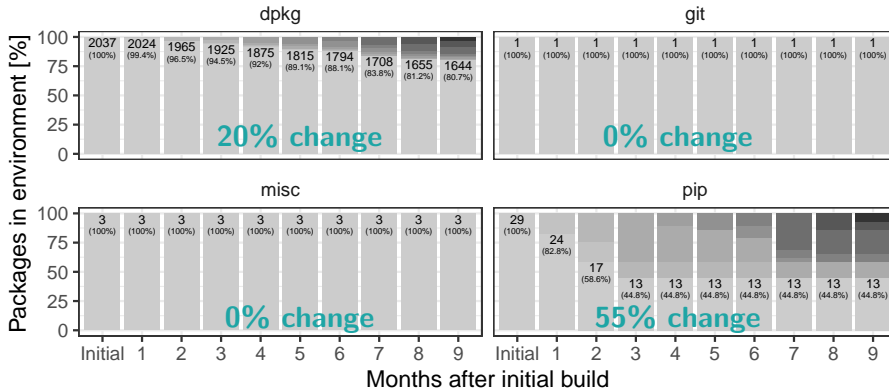
Preliminary Results – Per Artifact

Evolution of the packages versions over time for selected Dockerfile studied



Preliminary Results – Per Tool

Evolution of the packages versions over time

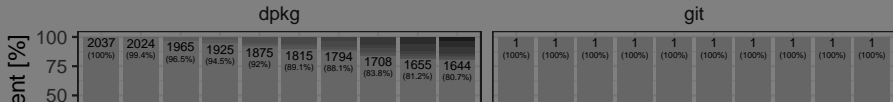


Month when the package version was introduced in the environment

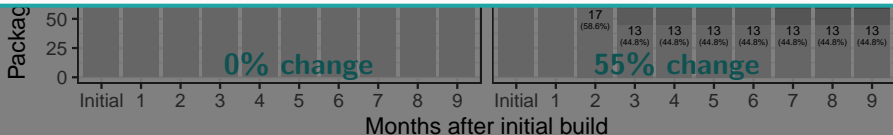


Preliminary Results – Per Tool

Evolution of the packages versions over time



Manual installations (git, misc):
more controlled \leadsto **more longevous**
but **challenging** to manage at scale (\leadsto Nix, Guix)



Month when the package version was introduced in the environment



How do the software environments produced by Dockerfiles from Artifacts evolve through time?

Preliminary Results

- **Software Env. changed within a month!** \leadsto **same period than the AE !**
- Only 5 artifacts 😞 (how significant / representative ?)

Future Work

- Design of the **large scale study** (How many artifacts? Which conferences?)
- Capture the hash of the base Docker image
- Other containerization tools? Other package managers?
- **Wanna help? Contact us!** 😊