

Popper Pitfalls: Experiences Following a Reproducibility Convention

Michael A. Sevilla

University of California, Santa Cruz
msevilla@soe.ucsc.edu

Carlos Maltzahn

University of California, Santa Cruz
carlosm@ucsc.edu

ABSTRACT

We describe the four publications we have tried to make reproducible and discuss how each paper has changed our workflows, practices, and collaboration policies. The fundamental insight is that paper artifacts must be made reproducible from the start of the project; artifacts are too difficult to make reproducible when the papers are (1) already published and (2) authored by researchers that are not thinking about reproducibility. In this paper, we present the best practices adopted by our research laboratory, which was sculpted by the pitfalls we have identified for the Popper convention. We conclude with a “call-to-arms” for the community focused on enhancing reproducibility initiatives for academic conferences, industry environments, and national laboratories. We hope that our experiences will shape a best practices guide for future reproducible papers.

CCS CONCEPTS

• **Software and its engineering** → **Reusability**; **Empirical software validation**; *Software evolution*; Collaboration in software development;

KEYWORDS

Software Testing ; Performance Engineering

ACM Reference Format:

Michael A. Sevilla and Carlos Maltzahn. 2018. Popper Pitfalls: Experiences Following a Reproducibility Convention. In *P-RECS’18: First International Workshop on Practical Reproducible Evaluation of Computer Systems, June 11, 2018, Tempe, AZ, USA*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3214239.3214243>

1 INTRODUCTION

The importance of reproducibility outweighs its difficulty. While wrestling with large amounts of data, millions of lines of code, and diverse developer groups, it is extremely hard to remember to make workflows understandable, let alone reproducible. But our research group has found that the benefits of reproducibility, from a research, educational, and productivity standpoint, make the pain points worth it. After all, the scientific method requires

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

P-RECS’18, June 11, 2018, Tempe, AZ, USA

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5861-3/18/06...\$15.00

<https://doi.org/10.1145/3214239.3214243>

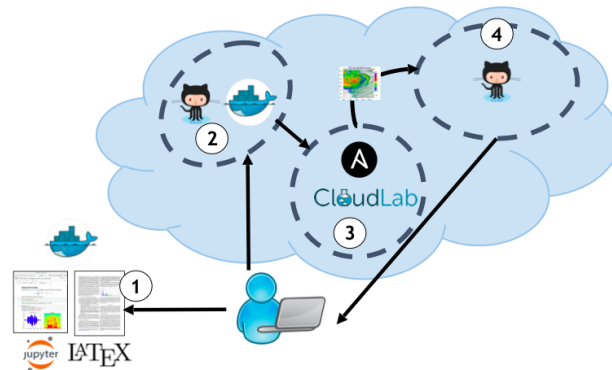


Figure 1: Illustration of the workflow used in our papers; Figure is adapted from [2]. For (1) - the visualization component - we use Jupyter, L^AT_EX, and Docker; for (2) - the code component - we use GitHub and Docker; for (3) - the multi-node component - we use Ansible and CloudLab; and for (4) - the data set component - we use GitHub.

such practices and other fields have been emphasizing the reproducibility of experimental results for *centuries*; it is about time that the field of computer systems change its *ad-hoc* workflows to be more reproducible.

Once we adopted this philosophy, the road to reproducible research and paper artifacts was not smooth. We use the Popper convention [2] because of its focus on open-source cluster management toolkits. Our workflow uses the DevOps tools shown in Figure 1 and is described in more detail in Section §2. We have produced four Popper-compliant papers, three of which have been published in the proceedings of top conferences.

While our experience with these state-of-the-art DevOps tools have been delightfully straightforward, which is no doubt a testament to the importance of reproducibility in real, large-scale clusters, our struggle composing these tools together in a coherent way has been a real challenge. From designing workflows with these tools, we have developed a list of pitfalls that wasted time and effort. Our contributions are as follows:

- the structure and workflow of some of our papers, using a baseline template designed for the Ceph [8] storage system (Section §2).
- three pitfalls for reproducibility that severely limited productivity in early papers; designing best practices that address these pitfalls before starting the project saved us time and effort in subsequent research explorations (Section §3).
- a call-to-arms for the downfall of practices that fly in the face of reproducibility. We highlight difficulties with academic

Paper [Reference]	Status	Reason
Malacology: A Programmable Storage System [6]	PASS	Reproducing results on different clusters not tested
Cudele: An API and Framework for Programmable Consistency and Durability in a Global Namespace [3]	GOLD	Tested on multiple clusters, READMEs for re-running experiments provided
Programmable Caches with a Data Management Language & Policy Engine [4]	FAIL	Security issues with public facing repositories and private networks/systems
Tintenfisch: File System Namespace Schemas and Generators [5]	GOLD	No performance results, just software replicability needed

Table 1: Papers following the Popper Convention. The statuses are from [1], where GOLD means results are reproducible, PASS means experiments run, and FAIL means experiment artifacts are available (but may not run). We graded the repositories based on our own internal evaluations – ideally, this evaluation would be done by someone outside the team, like a paper reviewer or a member of the conference committee. Over time, our Popper-compliance has improved, except for the work we did in [4], which faced security obstacles (see Section §4.2).

conferences, work at national laboratories, and work in industry because that is where we have the most experience (Section §4).

We give each of our paper repositories a status using the convention in [1]: GOLD means that results are fully reproducible, PASS means the experiments will run, and FAIL means that experiments are not guaranteed to execute. A FAIL status is still Popper-compliant because experiment artifacts and results are available. Over time, our Popper status has improved from PASS to GOLD because we were able to shape our best practices from our observed pitfalls. We also learned that going back and making published papers Popper-compliant is too difficult and is not incentivized, as described in Section §3.1. The one FAIL status we have can be attributed to security issues we had working at a national laboratory, as described in more detail in Section §4.2. It is our hope that the pitfalls we stumbled over can be used as a best practices guide for future articles.

2 POPPER-COMPLIANT PAPERS

We have produced four Popper-compliant papers, as shown in Table 1. These papers follow the reader/reviewer sample workflow outlined in [2] and shown in Figure 1. For the visualization component (1), we use Jupyter notebooks. The notebooks themselves are versioned with Git and users interact with local copies by cloning the repository and launching a Jupyter Docker container. The paper is written in \LaTeX and built with a Docker container. For the code component (2), both the source code for the system itself and the deploy/experiment code is stored on GitHub. When running experiments, we use Docker containers to isolate libraries and binaries. For the multi-node component (3), we use CloudLab machines and Ansible to script deployment and experiment orchestration. For the data set components (4), we use GitHub to store results files; our

inputs and results are small enough that we do not need a larger capacity. GitHub allows files up to 50MB and stores data on S3.

Our experiments start with a baseline. To describe the process, we reference our `ceph-popper-template`¹ set up on CloudLab. Users setup SSH keys and deploy CloudLab nodes using our CephFS Profile². The profile has the nodes automatically install Docker on bootup using our `install`³ script. After the nodes finish booting (*i.e.* their status on the CloudLab GUI is READY), users push SSH keys using a convenience script⁴.

The deploy code is based on `ceph-ansible`⁵, a tool that configures hardware and software for Ceph. We forked the project and made it less dependent on Python. To run an experiment, users log into the head node and clone the `ceph-popper-template` repository. This repository has submodules that point to `ceph-ansible` and our own custom roles; configuration files for our Ceph setup; and helper scripts written in bash that deploy Ceph and run the benchmarks. For more information on the Ceph template, see the README⁶ and for more information on the baseline and pipelines terminology, see the Popper Convention quickstart⁷.

An outline of this organization is shown in Figure 2, where solid lines are directory links and dashed lines indicate which directories a script uses. A sample experiment is in the `baseline/` directory, which was created using the Popper CLI. Users configure their cluster by specifying hostnames and login user names in the `hosts` file. Users also specify the Ceph services that should be deployed

¹<https://github.com/michaelsevilla/ceph-popper-template>

²<https://www.cloudlab.us/p/CephFS/CephFS-HEP>

³<https://github.com/michaelsevilla/ceph-popper-template/blob/master/hardware/cloudlab/install.sh>

⁴<https://raw.githubusercontent.com/michaelsevilla/ceph-popper-template/master/hardware/cloudlab/pushkeys.sh>

⁵<https://github.com/ceph/ceph-ansible/wiki>

⁶<https://github.com/michaelsevilla/ceph-popper-template>

⁷<http://falsifiable.us/>

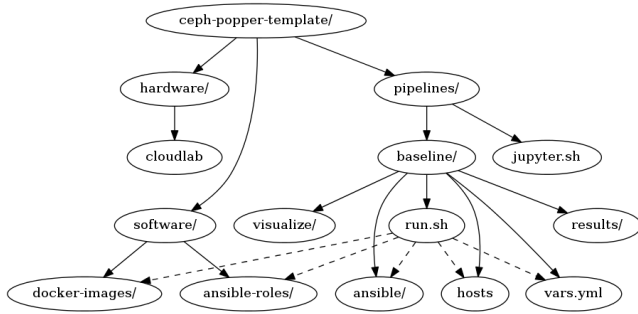


Figure 2: Organization of our experiment directories. `baseline/` is an example experiment and contains scripts populated with the Popper CLI. The `run.sh` script uses many directories and files to setup and run experiments, as indicated by the dashed lines.

using the `ansible/` directory⁸. This directory has code for deploying Ceph and its components, where the `*.yml` files are Ansible playbooks that start and configure components: `ceph.yml` starts Ceph and can be modified to specify which daemons to launch, `cleanup.yml` tears Ceph down, and `monitor.yml` starts daemons that monitor performance. High level configurations are in the `vars.yml` file. We separate these components into different playbooks so users can mix and match Ceph services. Similarly, the workloads directory has scripts for running the baseline benchmarks. The other files and directories are Ansible configuration files used by the playbooks. We have posted a tutorial on our blog⁹ to guide users through this setup. When complete, users execute the `run.sh` to start the job.

3 POPPER PITFALLS

To produce Popper-compliant papers, we needed to change our workflows and experimental procedures. We have shaped our best practices from the following pitfalls:

3.1 Reproducibility is not a 1st Class Citizen

Popper-compliance must be observed *throughout* the paper-writing process. Attempts to make published papers Popper-compliant failed for three reasons: (1) there is no incentive from the perspective of a researcher, (2) it is too hard to remember the experimental workflow, and (3) the artifacts cannot be added to the camera-ready article. We started trying to make [7] Popper-compliant¹⁰ but gave up because the effort would not have been rewarded, as the paper had already been presented and published.

Making reproducibility a 1st class citizen introduces a careful design decision: how does the user decide which results are exploratory and which are complete? We do not want to waste space by saving results that do not contribute to the paper but, at the same time, how do we know when an experiment is worthy of

being included in the paper? Discipline must be exercised throughout the paper-writing process to identify useful experiments and immediately polish them to be Popper-compliant.

Another pitfall we face is cross-cluster compatibility. After identifying results as being useful, we usually find cluster-specific files hard-coded throughout the tree, configuration files randomly propagated to different directories, visualization files reliant on data specified with hard-coded paths, and graphs that are not created automatically. To address this, we must exercise discipline to separate cluster-specific files from cluster-agnostic files; we do this with `configs_*/` directories in our `ceph-popper-template` repository.

Finally, organization and documentation must be required throughout the paper-writing process. We failed to do this in early papers and this causes the most work later on. System and experiment deploy code is rarely self-explanatory. At a minimum, each experiment directory must contain a README specifying how to run the jobs. To address this, we recommend relying heavily on Docker-Hub, GitHub, and collaboration; making things public incentivizes tidiness and organization. Using internal Docker or Git repositories hampers Popper-compliance and discourages community involvement.

3.2 Poorly-Defined Collaboration Roles

Collaboration speeds up code development and, as stated in the previous section, can be used as motivation for maintaining Popper-compliant repositories; but collaboration can also hamper the process. Researchers have different workflows and preferences. Committing system and experimental deploy code with different toolkits, documentation, and styles to the same repository can result in a confusing organizational structure. To address this, we recommend that the first author of the paper (1) produce a style guide, either by committing a couple of experiments or adding detailed READMEs and (2) act as the gatekeeper for all code. We recommend using GitHub because it nicely formats READMEs and supports pull requests, issue trackers, and project boards. Rather than having meetings to agree on experiment organization, we recommend iteratively combining code from collaborators with pull requests so the repository grows organically and quickly in the style approved by the first author.

Another obstacle to collaboration is conflicting updates, especially when editing the paper. We used to assign locks to people in an *ad-hoc* fashion over email but this is not scalable. Instead, we now use Git to manage conflicts. Before issuing a pull request, which notifies the first author, GitHub will notify committers of changes they made in parallel with other changes. Our policy is to force committers to resolve these merge conflicts before bothering the first author. While we still recommend issuing pull requests to make sure the first author approves of the style, the merge conflict workflow has been an effective way to avoid annoying the first author with destructive changes.

Git has not totally solved merge conflicts for papers. Papers have long sentences and sometimes a paragraph is written on a single line of text, without line breaks. Because Git was designed for code, which has many line breaks, it is hard to see changes when the document is structured as one long line. Adding line breaks partially solves the problem but adds new complexities. For example

⁸<https://github.com/michaelsevilla/ceph-popper-template/tree/master/pipelines/baseline/ansible>

⁹<http://programmability.us/mantle/blog2-ceph>

¹⁰<https://github.com/michaelsevilla/mds>

adding small amounts of text may cause lines to “spill” onto the next line, so a diff of the commits shows the entire paragraph as being changed. Similarly, users with different line break thresholds might re-arrange paragraphs differently, leading to the same problem. If Git or GitHub were to recognize that these source files were \LaTeX , the visual problems could be mitigated.

3.3 Failure to Maintain Pointers

Popper-compliant papers provide pointers to graph artifacts, including results, input files, and deploy code. But these links are problematic because they must be maintained over time, even years after the paper is published. We have had trouble maintaining these pointers because they (1) are ephemeral, (2) include pointers to other repositories, and (3) include pointers to different repository hosting sites. These problems arise because we are using GitHub and DockerHub as archival tools. Moving artifacts to a proper archival system (with long term storage and labeling semantics like DOIs) may be an option in the future, but GitHub works much better for research and exploration. Choosing between these options is a trade-off we have to make because there is no easy way for turning an exploration platform like GitHub into an archival system.

Ephemeral links and making sure that pointers are live is our biggest pain point. Small, seemingly benign changes end up confusing users wanting to reproduce our results. This is an artifact of how the internet was built, but making sure that these links are live is one of the hardest challenges for Popper-compliance. To combat this, we suggest maintaining continuous integration pipelines that read the paper source code and send REST requests to websites to ensure that 404 codes are not returned.

Pointers to other repositories can make the Popper-compliant repository confusing and difficult to navigate. For example, our Cudele paper had Git submodule pointers to 3 other GitHub repositories: one for Ansible deployment code for Ceph, one for Ansible deployment code for monitoring code written by our research lab, and one for our paper bibliography. Users unfamiliar with Ansible or \LaTeX would have trouble navigating these large code bases. Unfortunately, this practice is necessary to avoid duplicating functionality and ensuring future-proof modules. Our recommendation for addressing this problem is to minimize the number of repository pointers and to document them thoroughly, including the reasons that the pointer is necessary and what the other repository does.

Finally, pointers to different repository sites can make Popper-compliant papers difficult to understand. Again, this is a necessary evil because research systems can be large and complicated with many moving parts. For example, our Cudele paper uses a GitHub repository to maintain our source code for our modified Ceph version, a GitHub repository to maintain our paper and deploy code, a GitHub repository to maintain our common deploy code modules (for monitoring), a DockerHub repository for housing software images with compiled binaries for our modified Ceph version, and a CloudLab repository to house our base images. Again, the recommended solution for maintaining this complicated web of registry hubs is to thoroughly document the process.

4 COMMUNITY COOPERATION

We have shown and described reproducibility pitfalls, but equally important is community buy-in. Next, we outline practices that have been detrimental to our Popper-compliance initiatives.

4.1 Conference Requirements

The most common obstacle to our Popper-compliance efforts is double blinded submissions. The process is a burden, as we must create anonymous repositories and remove graph artifacts, but also blocks communication between researchers and reviewers. Providing evidence that the experiments work gives credence to the experiment and allows reviewers the chance to examine experiment parameters more closely. We go a step further and propose removing all anonymity from the review process, facilitating a communication channel between researchers and reviewers with the sole intent of improving the quality of the paper. We applaud the efforts of SC'18, IPDPS'18, and ASPLOS'18 as they move towards this approach with multiple revision rounds, but argue for more transparent review processes. We are encouraged that in the review of one of our papers, a reader specifically asked for source code and reproducibility artifacts; at that point, we gladly made Popper source links available.

A second obstacle is that many conferences lack a clear definition of reproducibility and replicability, which confuses both submitters and reviewers. One conference we submitted to had reviewers that posited that our paper reproducibility artifacts were out of scope while the submission website clearly had our definition of reproducibility. This confusion is frustrating and can lead to contentious reviews and rebuttals. To remedy the situation, we recommend emphasizing reproducibility initiatives to reviewers, even going as far as to reward papers that have clearly thought about reproducibility.

4.2 Industry/Laboratory Requirements

We understand the monetary incentives to propriety systems but working with code in these environments severely hampers our ability to make papers Popper-compliant. Some companies in industry keep all code repositories private. Furthermore, many companies have multiple repositories because development teams like using version control systems (*e.g.*, Git or SVN) and hosting services (*e.g.*, GitHub, GitLab, etc.) that they are familiar with. In larger companies that acquire startups, this is a big problem as every group of developers brings new ways of managing code. Obviously, this makes Popper-compliance impossible, except in a general sense.

We urge companies in industry to adopt a public, unified version control system and to provide reproducibility artifacts. Obviously, industry entities have very little incentive to do this. To incentivize this process for industry, we suggest that conferences provide recognition, either monetary or award-based, to companies that adopt this philosophy.

Another obstacle to Popper-compliance is security. Many systems in national laboratories require clearance and access is only granted to US citizens. In fact, many systems are not even connected to the internet to discourage contact with the outside world. While laboratories are expected to be research havens, their priorities are obviously security over open research practices. We recommend evangelizing reproducibility in these communities in

the hopes of appointing “reproducibility officers” that are internal to the national laboratories. These officers would have the security clearance to operate large HPC clusters and the expertise to verify the reproducibility of a paper’s artifacts and performance claims.

5 CONCLUSION

We have outlined the structure of our “reproducible” papers and showed the complexity of researching in this way. We iterated to a Popper-compliant paper-writing process after encountering numerous pitfalls, which we have documented and used to shape our best practices. While our process will never be perfect, we are encouraged with the improved speed and ease that our research progresses now that we have: (1) made reproducibility a 1st class citizen, (2) defined collaboration rules, and (3) made it a priority to maintain pointers. We hope that our call-to-arms for community cooperation effectively improves the state of reproducibility in the field.

Our future work is to quantify the efficiency improvements or degradations in efficiency due to reproducibility. An effective way to quantify productivity, in addition to the time spent on a problem, would help us more concretely identify when the reproducibility approach is useful or not. It might also help us explore other environments where reproducibility is useful, such as the classroom or in conference settings. Finally, it would highlight situations where the extra time spent setting up workflows and building paper artifacts is *worth it* for the ultimate payout of reproducibility when the paper is published.

ACKNOWLEDGMENTS

We thank the P-RECS reviewers for their helpful comments and suggestions. This work is supported by the Center for Research in Open Source Software (cross.ucsc.edu) and the NSF, under the award number 1450488.

REFERENCES

- [1] Ivo Jimenez, Sina Hamedian, Michael Sevilla, Noah Watkins, Carlos Maltzahn, Kathryn Mohror, Jay Lofstead, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. 2017. The Popper Experimentation Protocol. In *Series of Webinars on Reproducible Research*. <https://cross.ucsc.edu/wp-content/uploads/2017/09/jimenez.pdf>
- [2] Ivo Jimenez, Michael A. Sevilla, Noah Watkins, Carlos Maltzahn, Jay Lofstead, Kathryn Mohror, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. [n. d.]. The Popper Convention: Making Reproducible Systems Evaluation Practical. In *Proceedings of the International Parallel and Distributed Processing Symposium Workshop (IPDPSW '17)*.
- [3] Michael A Sevilla, Ivo Jimenez, Noah Watkins, Shel Finkelstein, Jeff LeFevre, Peter Alvaro, and Carlos Maltzahn. 2018. Cudele: An API and Framework for Programmable Consistency and Durability in a Global Namespace. In *Proceedings of the 32nd IEEE International Parallel and Distributed Processing Symposium (IPDPS '18)*.
- [4] Michael A. Sevilla, Carlos Maltzahn, Peter Alvaro, Reza Nasirigerdeh, Bradley W. Settlemyer, Danny Perez, David Rich, and Galen M. Shipman. 2018. Programmable Caches with a Data Management Language and Policy Engine. In *Proceedings of the International Symposium on Cluster, Cloud and Grid Computing (CCGrid '18)*.
- [5] Michael A. Sevilla, Reza Nasirigerdeh, Jeff LeFevre, Noah Watkins, Peter Alvaro, Margaret Lawson, and Jay Lofstead. 2018. *Tintenfisch: File System Namespace Schemas And Generators*. Technical Report UCSC-SOE-18-08. UC Santa Cruz.
- [6] Michael A. Sevilla, Noah Watkins, Ivo Jimenez, Peter Alvaro, Shel Finkelstein, Jeff LeFevre, and Carlos Maltzahn. 2017. Malacology: A Programmable Storage System. In *Proceedings of the European Conference on Computer Systems (EuroSys '17)*.
- [7] Michael A. Sevilla, Noah Watkins, Carlos Maltzahn, Ike Nassi, Scott A. Brandt, Sage A. Weil, Greg Farnum, and Sam Fineberg. [n. d.]. Mantle: A Programmable Metadata Load Balancer for the Ceph File System. In *Proceedings of the Conference on High Performance Computing, Networking, Storage and Analysis (SC '15)*.
- [8] Sage A. Weil, Scott A. Brandt, Ethan L. Miller, Darrell D. E. Long, and Carlos Maltzahn. [n. d.]. Ceph: A Scalable, High-Performance Distributed File System. In *Proceedings of the Symposium on Operating Systems Design and Implementation (OSDI '06)*.