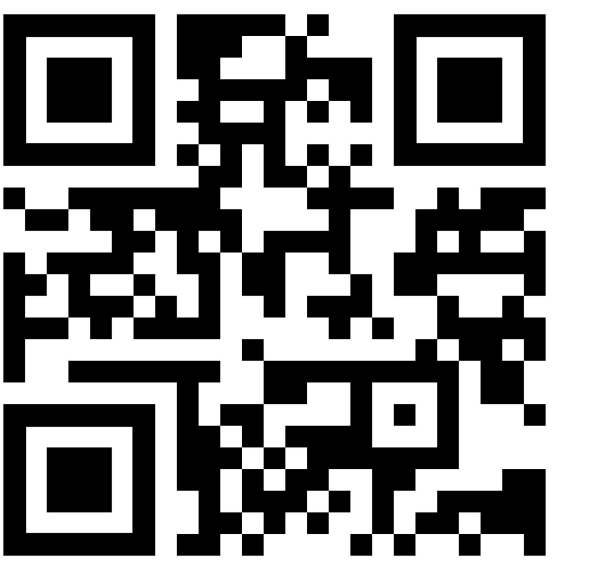




# Omnibenchmark for continuous and open benchmarking (in bioinformatics)



www.omnibenchmark.org

Daniel Incicau<sup>1</sup>, Reto Gerber<sup>1</sup>, Izaskun Mallona<sup>1,2</sup>, Almut Lütge<sup>1,2,3</sup>, Ben Carrillo<sup>1</sup>, Anthony Sonrel<sup>1,2</sup>, Charlotte Soneson<sup>2,4</sup>, and Mark D. Robinson<sup>1,2</sup>

1. Department of Molecular Life Sciences, University of Zurich | 2. SIB Swiss Institute of Bioinformatics

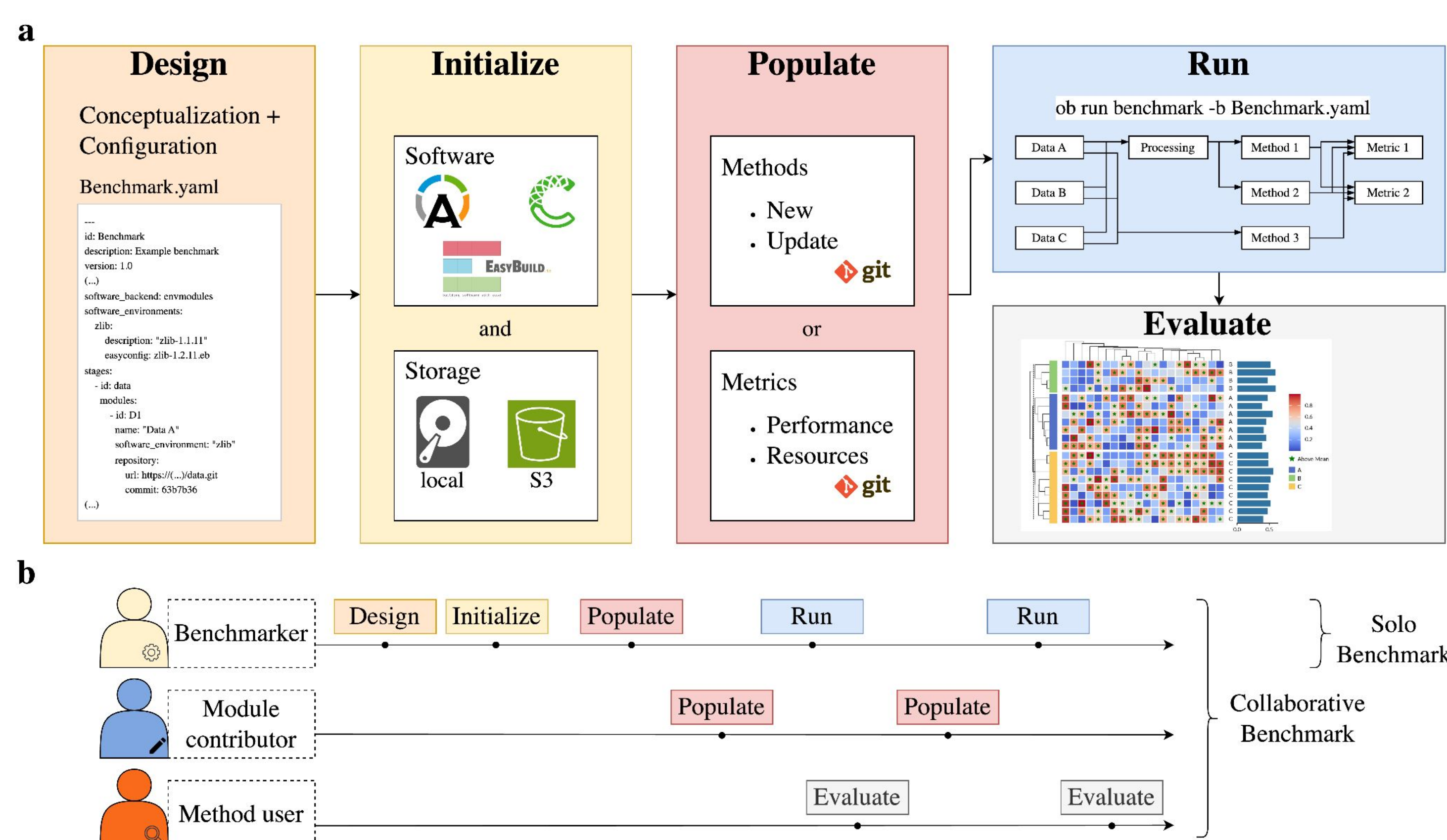
3. Swiss Data Science Center | 4. Friedrich Miescher Institute for Biomedical Research

## Introduction

**Omnibenchmark** provides **community-driven**, **extensible** and **continuously-updating** benchmarks. Omnibenchmark defines, executes and versions evaluation pipelines by leveraging a formal benchmark specification and a set of (reusable) benchmarking modules.

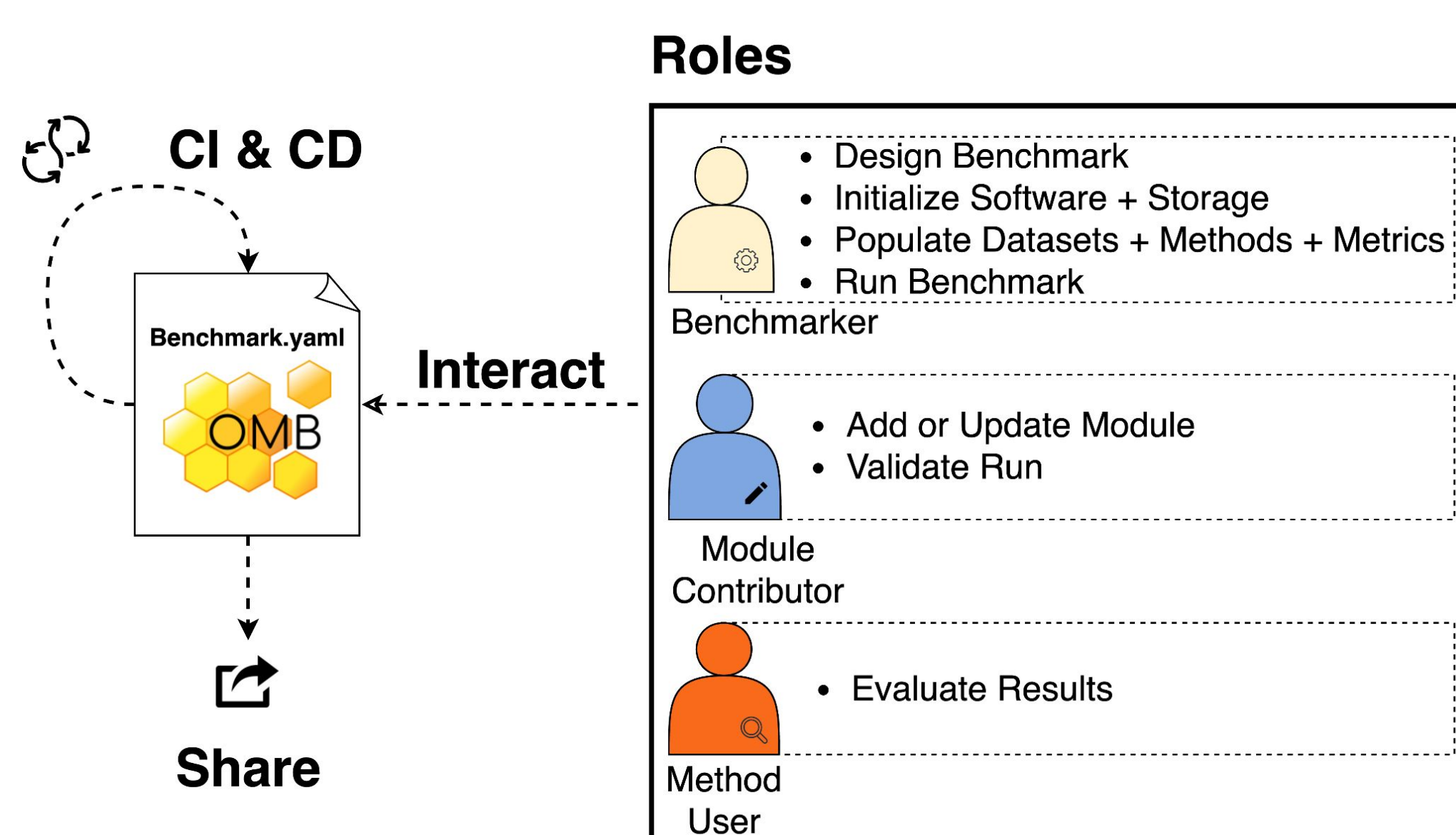
## Omnibenchmark Design and Usage

Figure 1a shows how to start a new benchmark using Omnibenchmark. After Conceptualization, the Configuration of the benchmark is specified in a single yaml file. Software and Storage need to be set up, followed by implementing methods and selecting metrics. A benchmark can be run by a single CLI command that generates the performance metrics. Finally, the results can be interactively visualized.



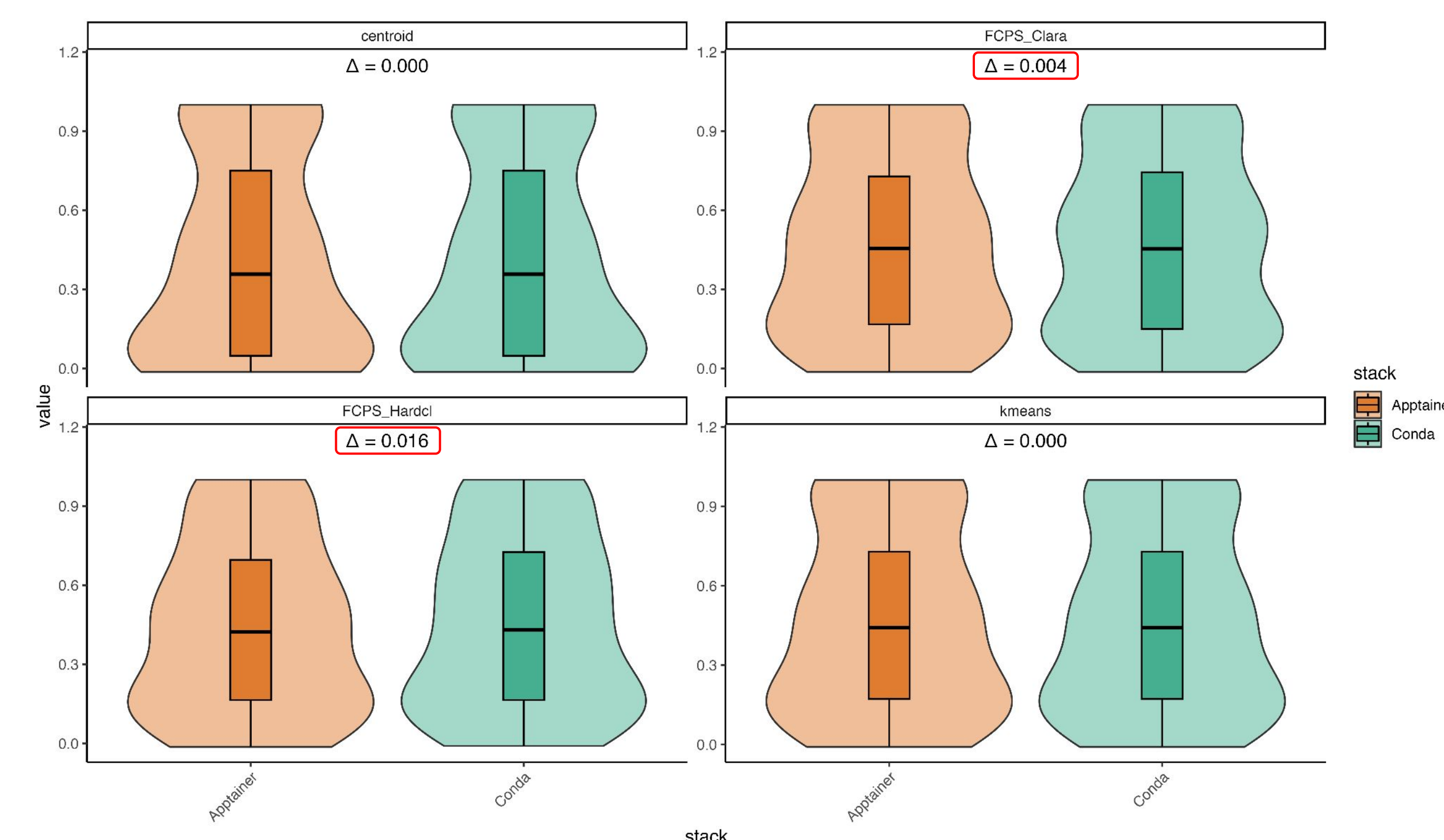
**Figure 1:** Omnibenchmark Workflow and Roles a) Benchmarking setup, from configuration to evaluation. b) Roles include Benchmarker (setup), Module Contributors (add/update methods and metrics), and Method Users (analyze results).

Different roles can be defined (Figure 1b), starting with the **'Benchmarker'**, who is responsible for starting and running a benchmark by initializing the configuration file, specifying the inputs and outputs of each stage, defining (initial) software stacks and providing storage. **'Module contributors'** can add or update datasets / methods / metrics to the benchmark and **'Method users'** evaluate the results to pick the best suited method for their use case.



**Figure 2:** Summary of Omnibenchmark Workflow and Roles

## Case study: clustering



**Figure 3:** Benchmark comparison of clustering methods executed in Conda and Apptainer environments. Violin plots display the distribution of Adjusted Rand Index (ARI) scores for each method, with pairwise median differences ( $\Delta$ ) shown above the plots. The comparison is performed across **62 datasets**, using the same code.

The supported software backends (**Apptainer**, **Conda**, **easybuild/Imod**) allow flexibility, but can result in slightly different results as shown for an example clustering benchmark (Figure 3).

In this figure, the choice of clustering methods highlights how software backends can impact clustering performance. While Centroid and K-means clustering yield consistent results across backends, the FCPS (Fundamental Clustering Problems Suite) shows variations depending on the backend used.

## Demo CI/CD

Easily trigger your own benchmark run by scanning the QR code below! This will initiate a predefined CI/CD workflow that uses the **Omnibenchmark CLI** on an example benchmark design.

### How it works:

1. Scan the QR code  
- Takes you to GitHub Actions
2. Ask the presenters to trigger the benchmarking pipeline  
- The workflows starts running the benchmark automatically
3. Monitor the results in real time  
- Track the execution and outcomes directly on Github



Try it out !

## Acknowledgements

- We thank the FOSS community for providing the building blocks of Omnibenchmark (Easybuild, git, Snakemake, conda, apptainer, etc)
- We thank the Renku team at the Swiss Data Science Center (SDSC)
- We thank members of Robinson Lab for their constructive feedback.
- This work has been supported with funding from the Swiss National Science Foundation (grants 200021\_212940 and 310030\_204869) as well as support from swissuniversities P5 Phase B funding (project 23-36\_14).