

Integration of Reproducible Methods into Community Cyberinfrastructure

David Tarboton^{1,1}, Tanu Malik^{2,2}, Jonathan Goodall^{3,3}, and Young-Don Choi^{3,3}

¹Utah State University

²DePaul University

³University of Virginia

November 30, 2022

Abstract

For science to reliably support new discoveries, its results must be reproducible. This has proven to be a challenge in many fields including fields that rely on computational methods as a means for supporting new discoveries. Reproducibility in these studies is particularly difficult because they require open, documented sharing of data and models and careful control of underlying hardware and software dependencies so that computational procedures executed by the original researcher are portable and can be run on different hardware or software and produce consistent results. Despite recent advances in making scientific work more findable, accessible, interoperable and reusable (FAIR), fundamental questions in the conduct of reproducible computational studies remain: Can published results be repeated in different computing environments? If yes, how similar are they to previous results? Can we further verify and build on the results by using additional data or changing computational methods? Can these changes be automatically and systematically tracked? This presentation will describe our EarthCube project to advance computational reproducibility and make it easier and more efficient for geoscientists to preserve, share, repeat and replicate scientific computations. Our approach is based on Sciunit software developed by prior EarthCube projects which encapsulates application dependencies composed of system binaries, code, data, environment and application provenance so that the resulting computational research object can be shared and re-executed on different platforms. We have deployed Sciunit within the HydroShare JupyterHub platform operated by the Consortium of Universities for the Advancement of Hydrologic Science Inc. (CUAHSI) for the hydrology research community and will present use cases that demonstrate how to preserve, share, repeat and replicate scientific results from the field of hydrologic modeling. While illustrated in the context of hydrology, the methods and tools developed as part of this project have the potential to be extended to other geoscience domains. They also have the potential to inform the reproducibility evaluation process as currently undertaken by journals and publishers.

Integration of Reproducible Methods into Community Cyberinfrastructure

Integration of Reproducible Methods into Community Cyberinfrastructure
 David G Tarboton (1)
 Tanu Malik (2)
 Jonathan L Goodall (3), Young-Don Choi (3)
 (1) Utah State University, (2) De Paul University, (3) University of Virginia

Problem

- Reproducible Crisis: Considerable research has documented difficulties in research reproducibility (Open Science, Stage et al., 2014)
- Not for science to reliably support new discoveries, its results must be reproducible
- Reproducibility of computational studies is particularly difficult because they require specific, documented settings of data and models and careful control of underlying hardware and software dependencies so that computational procedures executed by the original researcher are portable and can be run on different hardware or software and produce consistent results.

Goals of ReproBench EarthCube Project

- Advance computational reproducibility and make it easier and more efficient for geoscientists to generate, share, reuse and replicate scientific computations
- Advance the use of cloud software which can be applied across other disciplines composed of system libraries, code, data, environment and application processes so that the resulting computational research output can be shared and re-executed on different platforms
- Enable cloud native the HydroShare supported platform operated by the Consortium of Universities for the Advancement of Hydrologic Science in a Global Era for the hydrologic research community and demonstrate how to generate, share, reuse and replicate scientific results from the field of hydrologic modeling

An Actionable Approach to Reproducible Research

Computational Reproducibility

Computational Reproducibility requires establishing a progression from Repeatability through Reusability, Reproducibility, and Replicability, demanding increased time and effort.

Repeatability
 The same data and code are used to produce the same results.

Reusability
 The same data and code are used to produce different results.

Reproducibility
 The same data and code are used to produce the same results on different hardware or software.

Replicability
 The same data and code are used to produce the same results on different hardware or software.

The reproducibility taxonomy for complex computational studies (O'Leary et al., 2016)

Cyberinfrastructure Requirements

- Provision open data and results
- Provision Code
- Provision the computational environment

Repository and Compute

Other researchers in the context of hydrology, the pattern that links repository and compute capability and the methods and tools developed as part of this project have the potential to be extended to other geoscientific domains. They will have the potential to return the reproducibility-innovation process to currently practitioners by example and publishers.

Automatic Containerization of Execution Dependencies

http://hydroshare.com

Scientist Client obtains the requirements made to the host OS using a form of an requirement

1. C Container Scientist

2. Share Scientist

David G Tarboton (1)
 Tanu Malik (2)
 Jonathan L Goodall (3), Young-Don Choi (3)

(1) Utah State University, (2) De Paul University, (3) University of Virginia



PRESENTED AT:

EarthCube

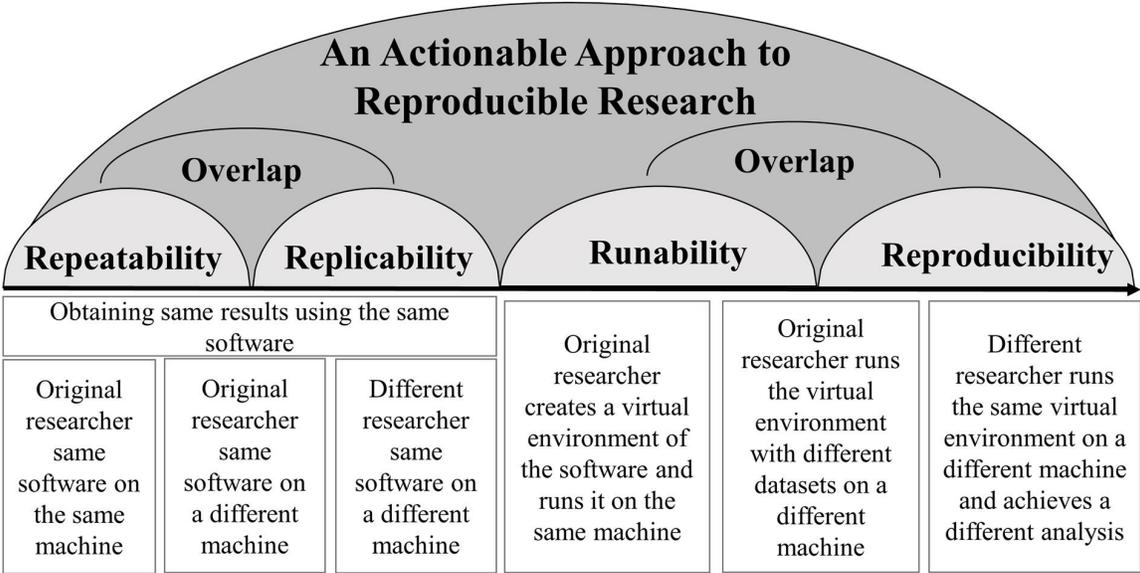
**2020 EarthCube Annual Meeting
 Virtual – June 18, 2020**

PROBLEM

- Reproducibility Crisis. Considerable research has documented difficulties in research reproducibility (Baker, 2016a,b; Stagge et al., 2019).
- Yet, for science to reliably support new discoveries, its results must be reproducible.
- Reproducibility of computational studies is particularly difficult because they require open, documented sharing of data and models and careful control of underlying hardware and software dependencies so that computational procedures executed by the original researcher are portable and can be run on different hardware or software and produce consistent results.

GOALS OF REPROBENCH EARTHCUBE PROJECT

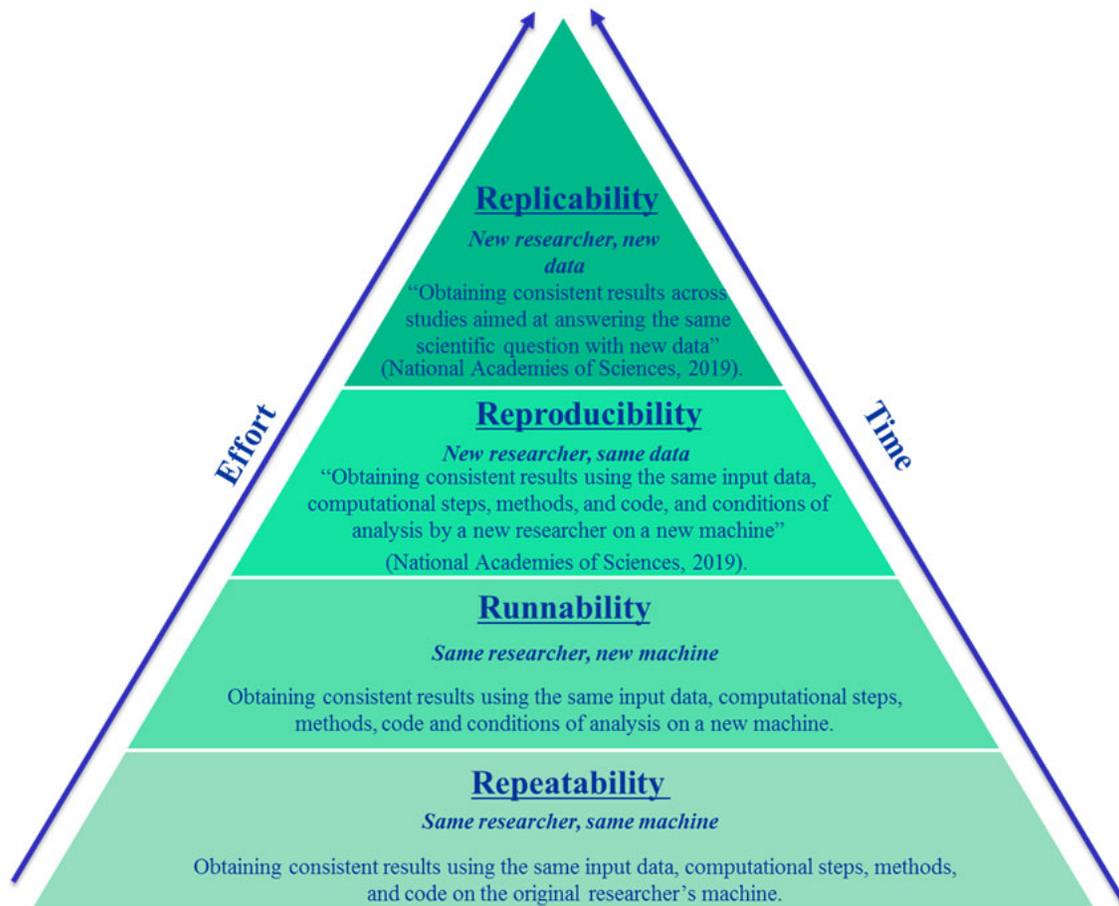
- Advance computational reproducibility and make it easier and more efficient for geoscientists to preserve, share, repeat and replicate scientific computations.
- Advance the use of Sciunit software which encapsulates application dependencies composed of system binaries, code, data, environment and application provenance so that the resulting computational research object can be shared and re-executed on different platforms.
- Deploy Sciunit within the HydroShare JupyterHub platform operated by the Consortium of Universities for the Advancement of Hydrologic Science Inc. (CUAHSI) for the hydrology research community and demonstrate how to preserve, share, repeat and replicate scientific results from the field of hydrologic modeling



Actionable approach for moving geoscience workflows from the Runability to Reproducibility level.

COMPUTATIONAL REPRODUCIBILITY

Computational Reproducibility requires establishing a progression from Repeatability, through Runnability, Reproducibility, and Replicability, demanding increased time and effort.



The reproducibility taxonomy for complex computational studies (Essawy et al., 2020).

Cyberinfrastructure Requirements

- Preserve input data and results
- Preserve Code
- Preserve the computational environment

Solution

- HydroShare platform for sharing and archiving data and models
- JupyterHub compute platforms linked to HydroShare for model execution
- Sciunit Software for encapsulating computational dependencies

HydroShare Data and Model Repository

- Manage data (and models and workflows) throughout research life cycle
- Share data, models, and other research products
- Permanent publication of data and models with citable digital object identifiers (DOIs)
- Fulfill Findable, Accessible, Interoperable, Reusable (FAIR) open data mandate

The screenshot displays the HydroShare website interface. At the top, the navigation bar includes 'HOME', 'MY RESOURCES', 'DISCOVER', 'COLLABORATE', 'APPS', and 'HELP', along with a 'SIGN IN' button. The main heading reads 'HydroShare is CUAHSI's online collaboration environment for sharing data, models, and code.' Below this is a 'Sign up now' button. A large banner image shows a stream with the word 'Discover' overlaid. A secondary navigation bar is visible above the resource page, which is titled 'TW Daniels Experimental Forest (TWDEF) Lidar'. The resource page includes metadata such as authors (Michaela Teich, David G. Tarboton), owners (Michaela Teich), resource type (Generic), storage (5.4 GB), creation date (Nov 17, 2016), last updated date (Nov 30, 2016), DOI (10.4211/hs.36f3314971a547bc8bc72dc60d9bd03c), and citation information. It also shows sharing status (Published), views (251), and downloads (11). An abstract section describes the lidar data collection and processing. A 'Resource Level Coverage' section provides spatial and temporal details, including a map of the study area in Utah. The 'Subject Keywords' section lists 'TW Daniels Experimental Forest', 'TWDEF', 'Lidar', 'DEM', and 'Snow Depth'.

CUAHSI and CyberGIS Jupyter for Water Gateways to computing

- Provide immediate value
 - What can I do now that I may not be able to easily do on my PC
- Model input data preparation
- Model execution
- Visualization and analysis (best of practice tools)
- Reduced needs for software installation and configuration (platform independence)
- Write and execute code in a Jupyter Notebook, acting on content of HydroShare resources and saving results back to HydroShare Repository
 - Collaboration
 - Access to enhanced computation (HPC, Big data)
- Enhanced trust in research through transparency, replicability and reproducibility

REPOSITORY AND COMPUTE

While illustrated in the context of hydrology, the pattern that links repository and compute capability and the methods and tools developed as part of this project have the potential to be extended to other geoscience domains. They also have the potential to inform the reproducibility evaluation process as currently undertaken by journals and publishers.

HydroShare

The screenshot shows the HydroShare website interface. At the top, there's a navigation bar with 'HYDROSHARE' and links for 'MY RESOURCES', 'DISCOVER', 'COLLABORATE', 'APPS', and 'HELP'. The main content area displays 'Hydrologic Terrain Analysis Jupyter Notebook' with metadata including authors (David Tarboton, Anthony Michael Catherine), creation date (June 03, 2018), and last update (June 05, 2018). Below the abstract, there's a 'Content' section showing a file named 'TauDEM.ipynb' (12.4 KB, Jupyter File). At the bottom, there are two stacked boxes: a green one labeled 'Django web framework' and a blue one labeled 'iRODS Network File System'.

JupyterHub

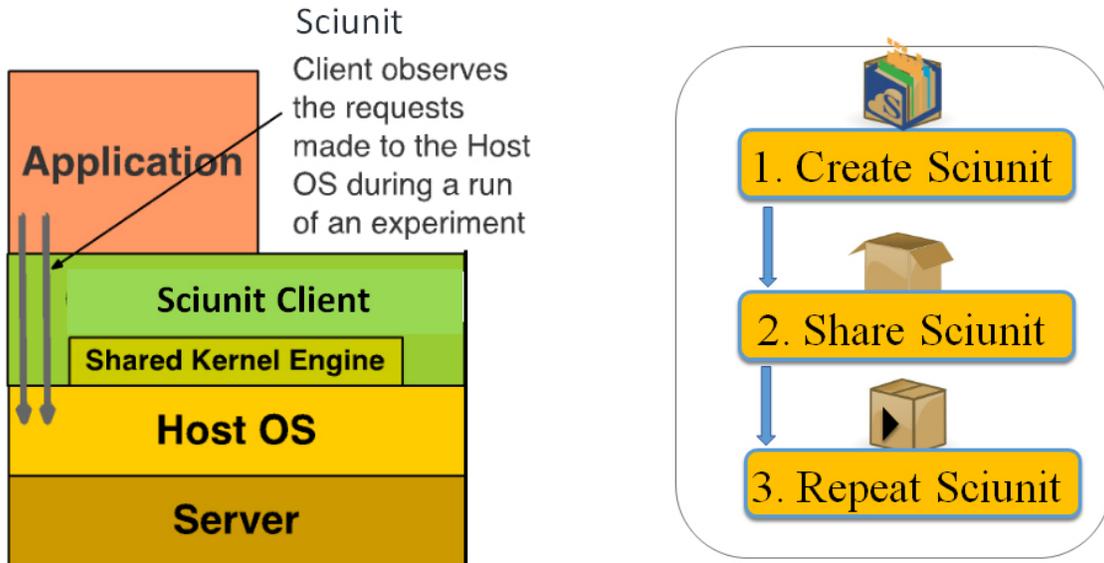
The screenshot shows a Jupyter Notebook interface titled 'TauDEM (autoexec)'. The notebook content includes an introduction to Hydrologic Terrain Analysis using TauDEM, followed by a list of tasks: '1- Preparation, libraries and getting oriented'. A yellow callout box contains the text: 'Write and execute code in a Jupyter Notebook, acting on content of HydroShare resources and saving results back to HydroShare Repository'. Below this, there's a code cell with Python code for finding files. At the bottom, there are two side-by-side plots: 'D8 Flow Direction' and 'D8 Slope'.

REST API

Oauth

AUTOMATIC CONTAINERIZATION OF EXECUTION DEPENDENCIES

<http://sciunit.run> (<http://sciunit.run>)



HydroShare Jupyter Notebook Resources that illustrate the use of Sciunit for reproducibility

- CHOI, Y. (2020). Sciunit SUMMA Result Reproduction Illustration, HydroShare, <http://www.hydroshare.org/resource/7d1403636fd3444c87e3c5b40b000b91> (<http://www.hydroshare.org/resource/7d1403636fd3444c87e3c5b40b000b91>) (This illustrates computational reproducibility using a model and computational environment encapsulated in a Sciunit stored in HydroShare. Details are described in Essawy et al., 2020)
- Choi, Y., J. Goodall, J. Sadler, A. M. Castronova, A. Bennett, T. Malik, B. Nijssen, Z. Li, S. Wang, M. Clark, D. Tarboton, M. Deeds (2020). EarthCube2020: An Approach for Open and Reproducible Environmental Modeling, HydroShare, <http://www.hydroshare.org/resource/75f31565dbd24c198450b9d37c6fcf74> (<http://www.hydroshare.org/resource/75f31565dbd24c198450b9d37c6fcf74>) (This illustrates the cycle involving the creation of a Sciunit container, saving to HydroShare and then re-execution of that container for computational reproducibility).

ABSTRACT

For science to reliably support new discoveries, its results must be reproducible. This has proven to be a challenge in many fields including fields that rely on computational methods as a means for supporting new discoveries. Reproducibility in these studies is particularly difficult because they require open, documented sharing of data and models and careful control of underlying hardware and software dependencies so that computational procedures executed by the original researcher are portable and can be run on different hardware or software and produce consistent results. Despite recent advances in making scientific work more findable, accessible, interoperable and reusable (FAIR), fundamental questions in the conduct of reproducible computational studies remain: Can published results be repeated in different computing environments? If yes, how similar are they to previous results? Can we further verify and build on the results by using additional data or changing computational methods? Can these changes be automatically and systematically tracked? This presentation will describe our EarthCube project to advance computational reproducibility and make it easier and more efficient for geoscientists to preserve, share, repeat and replicate scientific computations. Our approach is based on Sciunit software developed by prior EarthCube projects which encapsulates application dependencies composed of system binaries, code, data, environment and application provenance so that the resulting computational research object can be shared and re-executed on different platforms. We have deployed Sciunit within the HydroShare JupyterHub platform operated by the Consortium of Universities for the Advancement of Hydrologic Science Inc. (CUAHSI) for the hydrology research community and will present use cases that demonstrate how to preserve, share, repeat and replicate scientific results from the field of hydrologic modeling. While illustrated in the context of hydrology, the methods and tools developed as part of this project have the potential to be extended to other geoscience domains. They also have the potential to inform the reproducibility evaluation process as currently undertaken by journals and publishers.

REFERENCES

Baker, M., (2016a), "1500 scientists lift the lid on reproducibility," Nature News, 533(7604): 452-454, <https://doi.org/10.1038/533452a>.

Baker, M., (2016b), "Muddled meanings hamper efforts to fix reproducibility crisis," Nature News, <https://doi.org/10.1038/nature.2016.20076>

Essawy, B.T., Goodall, J.L., Voce, D., Morsy, M.M., Sadler, J.M., Choi, Y.D., Tarboton, D.G., Malik, T., 2020. A taxonomy for reproducible and replicable research in environmental modelling. *Environ. Model. Softw.* 104753. <https://doi.org/10.1016/j.envsoft.2020.104753>

Stagge, J. H., D. E. Rosenberg, A. M. Abdallah, H. Akbar, N. A. Attallah and R. James, (2019), "Assessing data availability and research reproducibility in hydrology and water resources," *Scientific Data*, 6: 190030, <https://doi.org/10.1038/sdata.2019.30>.