

Sharing and Preserving Computational Analyses for Posterity with *encapsulator*

Thomas Pasquier
University of Cambridge

**Matthew K. Lau and
Xueyuan Han**
Harvard University

**Elizabeth Fong and
Barbara S. Lerner**
Mount Holyoke College

**Emery R. Boose, Mercè
Crosas, Aaron M. Ellison,
and Margo Seltzer**
Harvard University

Editors: Lorena A. Barba,
labarba@gwu.edu;
George K. Thiruvathukal,
gkt@cs.luc.edu

Reproducibility has become a recurring topic of discussion in many scientific disciplines.¹ Although it might be expected that some studies will be difficult to reproduce, recent conversations highlight important aspects of the scientific endeavor that could be improved to facilitate reproducibility. Open data and open source software are two important parts of a concerted effort to achieve reproducibility.² However, multiple publications point out these approaches' shortcomings,^{3,4} such as the identification of dependencies, poor documentation of the installation processes, "code rot," failure to capture dynamic inputs, and technical barriers.

In prior work,⁵ we pointed out that open data and open source software alone are insufficient to ensure reproducibility, as they do not capture information about the computational execution, that is, the "process" and context that produced the results using the data and code. In keeping with the "open" culture, we defined open process as the practice of both sharing the source and the input data and providing a description of the entire computational

environment, including the software, libraries, and OS used for an analysis. We suggested the use of data provenance,⁶ formalized metadata representing the execution of a computational task and its context (for example, dependencies, specific data versions, and random or pseudorandom values), which can be captured during computation.

We view data provenance as key to addressing these issues, yet still insufficient. We need tools that leverage provenance to put capabilities, not complex metadata, into scientists' hands. We build on recent developments that address this need, such as executable papers⁷ and experiment packaging systems, for example, ReproZip.⁸ We propose a solution for scientists running small- to medium-scale computational experiments or analyses on commodity machines. Although tools exist to cover analyses done using spreadsheet programs (further discussed in the "Chal-

lenges” section), we intentionally do not cover that space, as it has inherent barriers to transparency and identification of the source of errors.^{9,10} Similarly, we do not attempt to address the reproducibility of large-scale computational analysis.

We present a “time capsule” for small- to medium-scale computational analysis. This time capsule is a self-contained environment that allows other scientists to explore the results of a published paper, reproduce them, or build upon them with minimal effort. We automatically curate the scientist’s code to extract only those elements pertinent to a particular figure, table, or dataset.

DATA PROVENANCE

Data provenance¹⁰ has the potential to address some of the challenges related to reproducibility. Indeed, to assess the validity or quality of information, it is necessary to understand the context of its creation. Unfortunately, digital artifacts frequently omit or hide much of the context in which they were created. As an example, many of us have been guilty of sharing code developed on our machines that our colleagues could not run; because we often work in the same environment for months or years, it’s easy to forget about software and libraries we have installed over time.

Meanwhile, small differences in a computational pipeline can lead to vastly different results. For example, different analyses of the same dataset of carbon flux in an Amazonian forest ecosystem differed in their estimates by up to 140 percent,¹¹ amounting to differences of up to 7 tons of carbon in an area the size of a football field. This example highlights the significant impact of small differences in code, especially when analyses or models contain user-defined or interactive (for example, multiplicative) terms. Seemingly small changes to inputs or in the computational pipeline can lead to large differences in results, impeding their reproducibility and verification.

Data provenance is a formal representation of the context and execution of a computation. This information is represented as a directed acyclic graph (DAG), a structure amenable to computational analysis. We use the World Wide Web Consortium (W3C) standard for data provenance: PROV-DM. Figure 1 shows a simple provenance graph. Vertices represent entities (representing data), activities (representing actions or transformations), and agents (representing users or organizations). In this figure, a process, controlled by Scientist Sarah, uses an executable function (a program) and an input file (data) to generate an output.

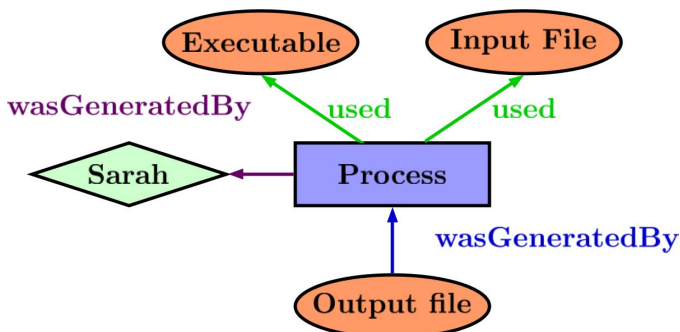


Figure 1. A simple W3C PROV-DM compliant provenance graph.

Provenance can be captured at various levels of a system, such as in libraries explicitly called by a program, in a language interpreter, in system libraries, or in the OS. The specific capture approach produces subtly different types of provenance: observed provenance is deduced by a system that monitors execution, whereas disclosed provenance is created explicitly by software that understands the semantics of the computations performed.¹² *encapsulator* uses observed provenance capture, which reveals the inner workings of an analysis script by collecting fine-grained provenance.

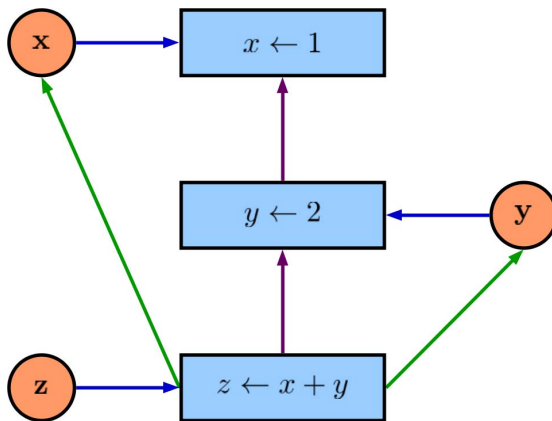


Figure 2. A simple provenance graph for an R script.

When provenance is captured for a scripting or programming language, the provenance DAG represents relationships among inputs, outputs, transient data objects, and statements. For example, Figure 2 illustrates a provenance graph of the R script shown in Listing 1.

```
x <- 1
y <- 2
z <- x + y
```

Listing 1. Simple R script example.

In Figure 2, the blue rectangles correspond to statements in the language; the orange circles correspond to data items (inputs, outputs, or transient objects); the purple arrows show the control flow, representing the precise sequence of steps taken while executing the program; and the blue and green arrows show data dependencies (the data used by an operation, and the data generated by an operation, respectively).

The provenance DAG illustrates data dependencies (what input generated a given output), software dependencies (on what libraries a script depends), and information about the structure of a program. We next discuss how we use provenance DAGs to generate a time capsule.

CREATING A TIME-CAPSULE

Provenance alone provides a “picture” of a computational context, yet we want to provide an active artifact that can reproduce a computational context: the time capsule. Figure 3 illustrates the two phases involved in creating a time capsule from the provenance collected during execution: curate the script to identify the precise lines of code and input data needed to produce a result, and build the time capsule containing the previously generated artifact and the environment necessary to reproduce it.

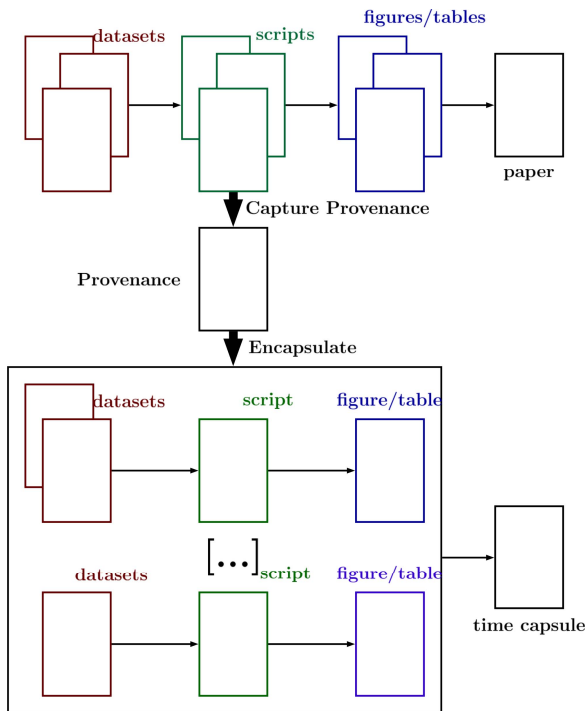


Figure 3. The encapsulation process.

Curating the Code

Science is, by its very nature, an iterative process. The task of cleaning and analyzing data is a stark example of this. The data obtained from scientific instruments or other measurements of the physical world are frequently a superset of the data a scientist wants to analyze. The first step in computation or analysis is often to “process” raw data to produce something that can be analyzed to answer a specific scientific question. This processing typically includes deciding how to handle missing data values, extracting parts of the data, computing new data from pieces of raw data, and so on. A scientist typically performs many such operations, not all of which end up being useful. Additionally, code evolves and accretes over time as scientists try different ways to interpret or analyze the data. False starts and abandoned analyses frequently persist in the final scripts that scientists use. The result is that code often contains a complex and evolving story of what transpired, rather than a clear, straight-line path from data to discovery. Although this history might be interesting, it might also lead to confusing and difficult-to-understand code.

The first phase of *encapsulator* takes as input the provenance of the computation’s execution, including all the false starts and abandoned attempts, and produces a curated script corresponding to the generation of a specific result. Such a curated script contains the minimum sufficient code to generate the output. Therefore, to understand a specific result, one can examine the curated version, rather than having to wade through potentially large amounts of irrelevant code.

To generate the minimal “cleaned” code, we analyze the provenance graph. Intuitively, the operations relevant to the generation of a figure or table are those connected in the DAG through data dependencies to the output. First, we trim the provenance graph by deleting control flow, considering only data dependencies. For example, the provenance graph illustrated in Figure 4 is transformed into the set of data dependency graphs shown in Figure 5.

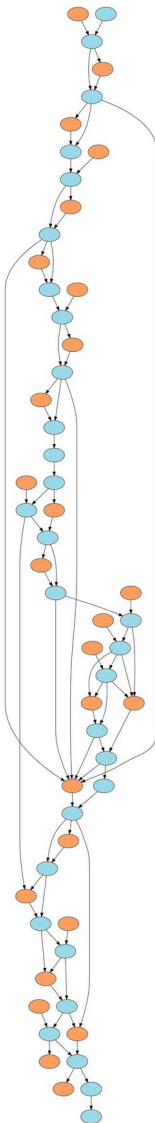


Figure 4. Provenance graph corresponding to a small R script (approximately 60 lines of code).

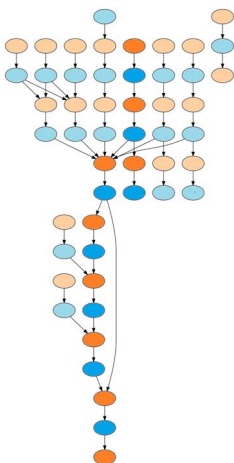


Figure 5. Data dependency transformation of the provenance graph shown in Figure 4.

In a data dependency graph, orange nodes represent inputs, outputs, or transient data, and blue nodes represent operations on data items. As we examine data dependencies in Figure 5, we alternate between data items and operations. The code necessary to generate an output (at the top of the figure) is the ordered set of operations present on all paths starting from the output in the original source code (more intense colored nodes in Figure 5 show an example of such a path for a given output). Similarly, the inputs necessary to generate an output are those encountered while traversing those paths. We generate the final, curated code by retaining all the operations on the paths in the graph leading to the output of interest, and then perform a final pass over the provenance DAG to identify all the required libraries. Once the final code has been generated, we run a source-code formatting tool (*formatR* for R scripts) to bring the code closer to best practices. We repeat these steps for every output of interest until we have generated a curated script for each. The inputs used to generate the selected outputs are identified and saved as part of the time capsule. We have made available (<http://provtools.org>) a standalone R library (*Rclean*; <https://cran.r-project.org/web/packages/Rclean>) implementing the mechanism described here.

Building the Time Capsule

Having shown how we produce curated scripts, we next explain how to construct a time capsule, leveraging freely available tools wherever possible. Our goal is to generate a self-contained environment that most scientists can use. This leads to the following requirements:

- The environment should present a user interface familiar to scientists.
- Encapsulation and use (de-encapsulation) of time capsules must require minimal technical expertise.
- The installation process itself must also require minimum intervention and technical knowledge.
- Time capsules, their installation, and re-execution must be platform-independent.

We demonstrate through a practical scenario how well we meet those requirements in the next section.

On the basis of those criteria, we selected virtual machines (VMs) as the self-contained environment for our time capsules (that is, their behavior and content is independent of the guest machine, and will remain identical over time). As one of the main barriers to reproducibility is technical, we want to avoid introducing additional technical complexity. Software such as VirtualBox (<https://www.virtualbox.org>) has made VMs an easy-to-use, “push button” technology, and it is possible to use a user-friendly interface to run a virtualized desktop with almost no technical knowledge. To most scientists, a VM will appear as a desktop environment similar to the one they use every day. To facilitate ease of adoption, we make sure that the time capsule contains all the tools scientists need to usefully interact with the computational process.

We use Vagrant (<https://www.vagrantup.com>) infrastructure and software to build, share, and distribute time capsules. Its VM provisioning is akin to that of Docker for containers. To provision a VM, one simply writes a script specifying the base VM (a preconfigured image), additional software, and files that should be installed. This is completely transparent to scientists: *encapsulator* generates a Vagrant file based on the information extracted from the provenance data in the previous phase. Although users can (optionally) customize the provision script, such customization should never be necessary. In the current prototype, the time capsule is Linux-based, as we leverage its package manager; other OSs present licensing challenges (discussed below). However, the creation of the time capsule itself can be done from experiments running on Windows, Mac, or any Linux distributions.

The provenance capture is achieved through program introspection using ProVR (<http://provtools.org/>). This presents some restrictions regarding the amount of system details that can be captured. In the current proof of concept implementation, we rely on the package manager of the Fedora Linux distribution (<https://fedoraproject.org/wiki/dnf>) to install the system dependencies required by a specific version of an R library. We are exploring the possibility

of complementing our provenance source using CamFlow (<http://camflow.org/>) to capture system-level provenance in the Linux OS. However, it must be noted that system-level provenance capture in closed source OSs remains a challenge.

During encapsulation, the scripts created in the first phase run in the time-capsule environment. Their outputs are compared to those from the original script (that is, the one run on the host machine) to ensure that they are identical. Once the encapsulation is finished, the VM is packaged, ready to be shared. This VM contains individual R scripts for each selected figure, along with the datasets used as inputs. The current prototype relies on Vagrant’s cloud platform to host the VM.

USING *ENCAPSULATOR*

Consider the following scenario: Sarah is a young and brilliant scientist who would like to make her research results available to the community, allow reviewers to easily verify her results, and encourage others to build on them. John is a scientist from a near future who wishes to use Sarah’s results. Professor O’Brien is a reviewer, interested in verifying Sarah’s findings.

The “messycode” examples (<https://github.com/ProvTools/encapsulator>) illustrate several “lazy coding practices” that scientists, including Sarah, often use when writing code for models and analyses:

- near stream-of-consciousness coding that follows a train of thought in script development,
- output to console that is not written to disk,
- intermediate objects that are abandoned,
- library and new data calls throughout the script,
- output written to disk but not used in final documents,
- code that is not modularized, and
- code that is syntactically correct but not particularly comprehensible.

At this stage, we assume that Sarah has finished her computations, built the figures and tables for her paper, and has the paper ready for submission. She is aware of research data repositories such as Dataverse repositories (<https://dataverse.org/>), and source-code repositories such as GitHub, but she knows that they might not be sufficient to make her code truly reusable. In the past, when she tried to reuse code written by other scientists, she often discovered that it was poorly documented and hard to use. She also constantly found herself baffled by questions such as what external packages the computation depends on, where to obtain those dependent files and libraries, and what parameters were used to obtain the published results. Trying to figure out these details resulted in her wasting countless hours. She would like to save other scientists from these challenges, so that they can more easily build upon her work.

Sarah wants a picture of the context of her computations that allows anyone to reproduce them. Provenance captured by tools such as *provR* (<http://provtools.org/>) for R scripts contains the following information, represented as nodes or node attributes in a DAG:

- inputs,
- outputs,
- transient data objects and their values,
- operations, and
- library dependencies.

This information facilitates depiction of the development environment, accurately capturing, for example, random seeds used and the version of a library that was required by the system. Although this picture is important, it could prove difficult for John or Professor O’Brien to use it to create an environment in which Sarah’s computations can be reproduced. They might not have the required expertise, or the required version of a library might have become unavailable. Thus, Sarah wants her experiments to be preserved in a time capsule.

Sarah decides to use *encapsulator*. She needs to install it and its dependencies: VirtualBox and Vagrant. On her Mac laptop, she can do this:

```
brew install ruby
gem install encapsulator
encapsulator --install mac
```

Listing 2. Installing *encapsulator* and its dependencies.

The next step is to examine her R script and determine what outputs she wants to include in her time capsule. She can find out what the possibilities are using *encapsulator*'s info capability:

```
encapsulator --info sarah.R
```

Listing 3. Obtaining a summary of an R script.

This generates the following output:

```
Files
-----
Input july_biomass_survey.csv
Input dataset_v2_june_from_collaborator1.csv
Output save1.csv
Output fig1_biplot.png
Output fig1_biplot_v2.png
Output fig2_biplot.png

Packages
-----
base v3.4.0
gdata v2.18.0
lattice v0.20-35
permute v0.9-4
txtplot v1.0-3
vegan v2.4-3
```

Listing 4. Example of an R script summary.

Sarah included only `fig1_biplot_v2.png` and `fig2_biplot.png` in her article, so she wants to generate a time capsule containing only the code (<https://github.com/ProvTools/encapsulator>) needed to generate those two images:

```
encapsulator --encapsulate sarah/experiment sarah.R fig1_biplot_v2.png fig2_bip-
lot.png
```

Listing 5. Creating the time capsule.

Once *encapsulator* has finished building the time capsule, all that is left for Sarah to do is upload it to her Vagrant cloud account.

A few months later, Professor O'Brien is reviewing Sarah's paper and wants to understand her analysis. He sees that Sarah has used *encapsulator* to share her work. As Sarah did in her workflow to produce the published results, he can easily install it on a Linux machine:

```
sudo apt install ruby
gem install encapsulator
encapsulator --install ubuntu
```

Listing 6. Installing *encapsulator* and its dependencies.

Once it is installed, he retrieves Sarah's work by running:

```
encapsulator --decapsulate sarah/experiment
```

Listing 7. De-encapsulating a shared environment.

encapsulator manages the VM download and start-up transparently. After a short time, a window appears on Professor O'Brien's desktop presenting him with the virtual desktop shown in Figure 6. In this environment, he has access to familiar tools and can work without difficulty. Further, the code that he examines for each figure is about a dozen lines of clean code, not Sarah's original 60 lines of messy code. Naturally, *encapsulator* can handle longer and more complex scripts.

John reads Sarah's article five years after its publication. Using the same sequence of commands that Professor O'Brien did, he is able to get the time capsule running on his laptop, and the environment in the VM is identical to what it was at the time of publication. John can get to work easily without worrying about the problem of outdated dependencies (such as old library versions that are no longer available for download).

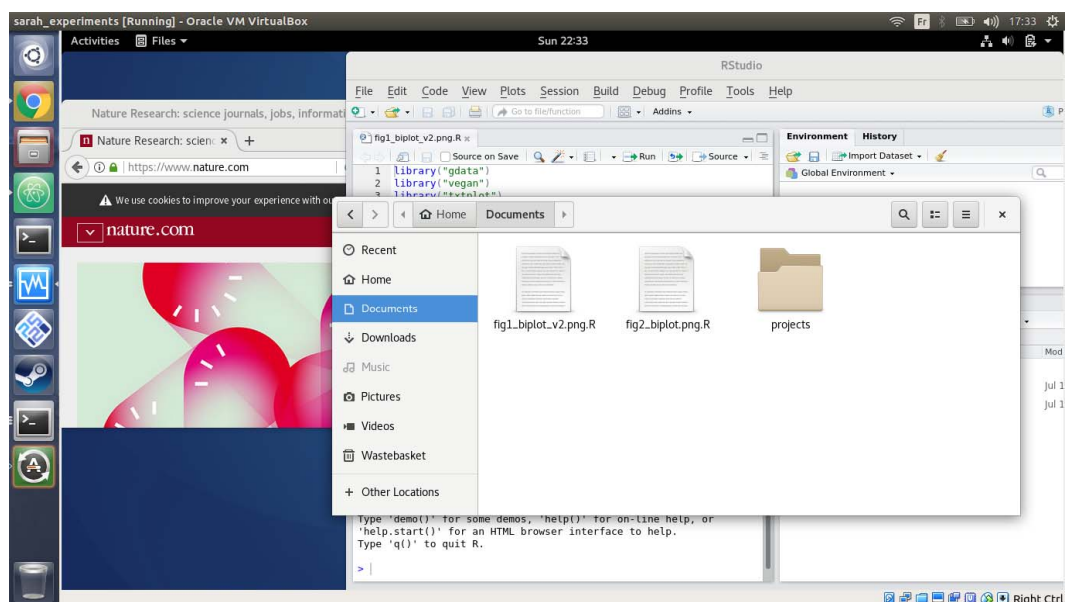


Figure 6. The time capsule running on Professor O'Brien's machine.

CHALLENGES

Domain-Specific Environment

Our time capsule comes with a generic environment, including some tools generally used for data analysis to provide an easy-to-use, familiar interface. In future versions, based on domain-scientist feedback, we will provide platforms containing standard toolsets specific to domains (“ecology”, “genetics”, “chemistry,” and so forth).

Time-Capsule OS

The current version of *encapsulator* uses Linux, in particular the package management system, to build a time capsule. Although a large number of tools used by scientists are available on Windows, Mac, and Linux, some tools might be available only on specific platforms. Furthermore, distributing Mac and Windows capsules introduces licensing issues (proprietary software in research is a complex topic¹³). At this stage, one can build a capsule on any platform, but the capsule itself is Linux-based. This might not pose a major obstacle for domain scientists whose analytical workflows occur almost entirely within an integrated development environment (IDE) such as RStudio, since these IDEs are supported on all major OSs and appear nearly identical across platforms.

Language Support

Our current prototype supports only the R programming language. We intend to incorporate support for additional languages used in data analysis, including Python and provenance capture libraries such as noWorkflow.¹⁴ Because *encapsulator* uses the PROV-JSON standard for data provenance, any provenance capture tool with a statement-level granularity for any language could be used to generate a capsule. Furthermore, it should be possible to support individual workflows that use multiple languages, which are becoming more common in some domains.

Integration with IDEs

Although they are relatively simple to use, command-line interfaces are daunting to some users. We are investigating integrating *encapsulator* into existing, commonly used IDEs, such as an *encapsulator* add-in for RStudio, a common IDE for R (<https://rstudio.github.io/rstudioaddins/>). Many researchers use spreadsheet programs for their data management and analysis. Although the feasibility and sufficiency of capturing provenance for such workflows has been demonstrated,¹⁵ and encapsulation is therefore also theoretically possible, we argue that these methods are inherently unstable since they typically rely on proprietary software with complex underlying data structures. Additionally, best practices for data science typically conflict with spreadsheet-based workflows that tend to lead to informal, and often inaccurate, data management and analysis.

Out-of-Tree Libraries

Many obscure libraries might not be available through the package management system, either a specific Linux distribution or a programming language, such as CRAN (<https://cran.r-project.org>) for R packages. We are investigating ways to handle such library dependencies. Those that do not have dependencies are relatively easy to handle by building and installing the package during the encapsulation process. Others that use alternative package managers, such as Bioconductor (<https://www.bioconductor.org>), are also relatively easy to handle. However, those with complex third-party dependencies without formal definitions are more difficult to support.

Nondeterministic Processes

Some scripts use pseudorandom number generators and two runs might not produce identical results. We plan to incorporate the ability to reproduce such results in a future release once the provenance capture system records random values; however, a more serious issue is nondeterminism introduced by concurrency. This could be ameliorated during the curation phase by producing scripts that enforce ordering. It might be preferable to enhance how we assess whether a given result produced within the time capsule is correct. In the current proof of concept, the results must be identical to those produced on the host machine. However, it might be reasonable to verify that the results meet some statistical property instead (for example, within δ of the original results). We recognize that this is not a trivial task and that significant investigation is required to determine a suitable path forward.

Long-Term Archival

There are two major assumptions that *encapsulator* makes about availability of a time capsule for long-term archival: the continued existence of the Vagrant cloud, and x86-64 virtualization. The first issue can be addressed by replicating the time capsule in a trusted archival repository. One option that we plan to explore in future work is to publish the time capsule in a Dataverse repository as a “replication dataset,” automatically assigning a DOI and minimal citation metadata and generating a formal persistent data citation for the time capsule. The second issue is more complex, so the answer is speculative. Virtualization depends on the remaining life span of the x86-64 architecture and whether the concerned time capsule will have any relevance after that. This last point is interesting to ponder, as preservation of our digital world is an issue¹⁶ that goes beyond science and reproducibility. Artifacts of our modern culture are already disappearing (such as video games and digital publications), which is an important sociocultural issue beyond the scope of our current project.

Container Support

Although we claim that tools such as Docker are not ideal to reduce the technical barriers to reproducibility for scientists, they are useful for automating the repetition of results. As Vagrant supports container provisioning, *encapsulator* can handle such targets. However, one should also remember that while containers are lighter, they are not as self-contained as VMs. Indeed, containers run over the kernel of their host machine; if change to the kernel were to affect results, then reproducibility could not be guaranteed.

CONCLUSION

We introduce *encapsulator*, a sophisticated yet simple toolbox that uses the provenance of computational data analysis to produce a time capsule in which computational workflows can be re-run and modified. This tool is designed to require minimal overhead for integration into a user’s workflow and limited technical expertise. When viewed within the context of increasing computational demands of all disciplines, *encapsulator* provides a key tool for facilitating transparent research at a crucial time for science.

SOFTWARE ENGINEERING PRACTICES

All software presented in this column is open source under GPL v3, and available at <http://provtools.org>, or directly through GitHub (<https://github.com/ProvTools>). The latest version (at the time of submission) can be referenced with doi:10.5281/zenodo.1199232, and is distributed via the RubyGems service (<https://rubygems.org/gems/encapsulator>). The software presented here remains under development and is subject to change. Matthew K. Lau should be contacted for any additional information about the ProvTools ecosystem. Further details about continuous integration and engineering practices are available in the README.md files of the individual components.

SIDEBAR: ALTERNATIVES TO ENCAPSULATOR

Some systems are designed to reproduce complex workflows running on grid or cloud infrastructures (for example, Kepler¹⁷), and fill a related, but distinct niche. Indeed, *encapsulator* is intended to support research run on single commodity machines, which accounts for a significant proportion of research results in a number of fields. Systems designed for particular domains already exist (for example, GenePattern¹⁸ and Galaxy¹⁹), but the role of *encapsulator* is to provide a general approach.

ReproZip⁸ and CDE²⁰ are directly comparable to *encapsulator*. However, they use system calls to identify dependencies and package experiments. Therefore, computations must first be run in Linux before they can be packaged. This might prove problematic for many scientists who do not use Linux. *encapsulator* relies on language-level observed provenance and is not subject to such limitations.

The main difference between *encapsulator* and alternative tools is its ease of use. Modifying packaged computations generated by the alternatives might require a relatively high level of technical skill. *encapsulator* builds a fully functional, self-contained environment that is easy for scientists to navigate. The list presented here is succinct, but we maintain online a list of open source provenance tools, including some designed for reproducibility and replication purposes (<https://projects.iq.harvard.edu/provenance-at-harvard/tools>).

ACKNOWLEDGMENTS

This work was supported by the NSF (grant no. SSI-1450277 End-to-End Provenance, and grant no. ACI-1448123 Citation++). More details about those projects are available at <https://projects.iq.harvard.edu/provenance-at-harvard>.

The peer reviewers for this manuscript were Professor Lorena Barba (School of Engineering and Applied Science, George Washington University) and Professor Carl Boettiger (Department of Environmental Science, Policy and Management, University of California, Berkeley). Both helped to clarify the terminology used around reproducibility, and Professor Boettiger helped us clarify the extent of the provenance captured.

REFERENCES

1. M. Baker, "1,500 Scientists Lift the Lid on Reproducibility," *Nature*, vol. 533, no. 7604, Nature, 2016, pp. 452–454.
2. J. Gezelter, "Open Source and Open Data Should Be Standard Practices," *J. of Physical Chemistry*, ACS, 2015; doi.org/0.1021/acs.jpcclett.5b00285.
3. D Garijo et al., "Quantifying Reproducibility in Computational Biology: The Case of the Tuberculosis Drugome," *PloS one*, vol. 8, no. 11, PloS one, 2013, p. e80278.
4. L. Joppa et al., "Troubling Trends in Scientific Software Use," *Science*, vol. 340, no. 6134, American Association for the Advancement of Science, 2013, pp. 814–815.
5. T. Pasquier et al., "If These Data Could Talk," *Nature Scientific Data*, Nature, 2017; doi.org/10.1038/sdata.2017.114.
6. L. Carata et al., "A Primer on Provenance," *Comm. ACM*, vol. 57, no. 5, ACM, 2014, pp. 52–60.
7. R. Strijkers et al., "Toward Executable Scientific Publications," *Procedia Computer Science*, vol. 4, Elsevier, 2011, pp. 707–715.
8. F. Chirigati et al., "Reprozip: Computational Reproducibility with Ease," *Proc. Int'l Conf. Management of Data*, 2016, pp. 2085–2088.
9. J. Cunha et al., "Towards a Catalog of Spreadsheet Smells," *Proc. Int'l Conf. Computational Science and Its Applications*, 2012, pp. 202–216.
10. M. Ziemann, Y. Eren, and Q. El-Osta, "Gene Name Errors are Widespread in the Scientific Literature," *Genome Biology*, vol. 17, no. 1, 2016, p. 177.

11. A. Ellison et al., “An Analytic Web to Support the Analysis and Synthesis of Ecological Data,” *Ecology*, vol. 87, no. 6, 2006, pp. 1345–1358.
12. U. Braun et al., “Issues in Automatic Provenance Collection,” *Proc. Int'l Conf. Provenance and Annotation of Data (PAW 06)*, 2006, pp. 171–183.
13. A. Gambardella and B. Hall, “Proprietary Versus Public Domain Licensing of Software And Research Products,” *Research Policy*, vol. 35, no. 6, 2006, pp. 875–892.
14. L. Murta et al., “noWorkflow: Capturing and Analyzing Provenance of Scripts,” *Proc. Int'l Provenance and Annotation Workshop*, 2014, pp. 71–83.
15. H. Asuncion, “In situ Data Provenance Capture in Spreadsheets,” *Proc. IEEE Int'l Conf. eScience*, 2011, pp. 240–247.
16. K. Lee et al., “The State of the Art and Practice in Digital Preservation,” *J. Research Nat'l Institute of standards and technology*, vol. 107, no. 1, 2002, p. 93.
17. I. Altintas et al., “Kepler: An Extensible System for Design and Execution of Scientific Workflows,” *Proc. 16th IEEE Int'l Conf. Scientific and Statistical Database Management*, 2004, pp. 423–424.
18. M. Reich et al., “Genepattern 2.0,” *Nature Genetics*, vol. 38, no. 5, 2006, pp. 500–501.
19. B. Giardine et al., “Galaxy: A Platform for Interactive Large-Scale Genome Analysis,” *Genome Research*, vol. 15, no. 10, 2005, pp. 1451–1455.
20. B. Howe, “CDE: A Tool for Creating Portable Experimental Software Packages,” *Computing in Science & Engineering*, vol. 14, no. 4, 2012, pp. 32–35.

ABOUT THE AUTHORS

Thomas Pasquier is a research associate at the University of Cambridge’s Department of Computer Science, a research fellow at St. Edmund’s College (University of Cambridge) and an associate of Harvard University’s Center for Research on Computation and Society. His research interests include the design of more accountable and transparent computer systems. Pasquier received a PhD in computer science from the University of Cambridge. He is a member of IEEE and ACM. Contact him at tfjmp@seas.harvard.edu.

Matthew K. Lau is a postdoctoral research fellow at Harvard University’s Harvard Forest Long-Term Ecological Research Site, and a member of the Ecological Society of America. His research interests include the ecological and evolutionary dynamics of ecological networks. A long-time supporter of open source programming tools for scientific analyses, he has taught undergraduate- and graduate-level courses and workshops in statistical programming. Lau received a PhD in biology from Northern Arizona University. Contact him at matthewklau@fas.harvard.edu.

Xueyuan Han is a second-year PhD student in computer science at Harvard University. His research interests include applying data provenance analysis in various contexts, including data management, cybersecurity, and machine learning. Han received a BS in computer science, as well as an Outstanding Bachelor of Science Department Award, from the University of California, Los Angeles. Contact him at hanx@g.harvard.edu.

Elizabeth Fong is a software developer and researcher at Mount Holyoke College. In between graduating from and working at Mount Holyoke College, she was a data engineering fellow at Insight. Her research interests include data provenance and data engineering, as well as their applications in fields such as ecology, medicine, and biochemistry. Fong received a BA in computer science with a minor in molecular biology from Mount Holyoke College. Contact her at fong22e@mholyoke.edu.

Barbara S. Lerner is an associate professor of computer science at Mount Holyoke College. Her research interests include automated collection of data provenance to support reproducibility, program understanding, and debugging. She is also interested in the development of software/hardware solutions to help people with visual impairments. Lerner received a PhD in computer science from Carnegie Mellon. Contact her at blerner@mholyoke.edu.

Emery R. Boose is an information manager and senior scientist at the Harvard Forest Long-Term Ecological Research Site. His research interests include data provenance, ecoinformatics, hurricane modeling, meteorology, and hydrology. Boose received a PhD in Sanskrit and Indian studies from Harvard University. Contact him at boose@fas.harvard.edu.

Mercè Crosas is the chief data science and technology officer at the Institute for Quantitative Social Science (IQSS) at Harvard University. She has more than 10 years of experience leading the Dataverse project, and more than 15 years of experience building data management and analysis systems in academia and biotechnology companies. She is part of numerous committees and working groups focused on data sharing, research data management, data citation, and data standards. Crosas received a PhD in astrophysics from Rice University. Contact her at mcrosas@iq.harvard.edu.

Aaron M. Ellison is the senior research fellow in ecology in Harvard's Department of Organismic and Evolutionary Biology and a senior ecologist at the Harvard Forest Long-Term Ecological Research Site. His research focuses on the assembly, disassembly, and reassembly of ecosystems following natural and anthropogenic disturbances, and on the relationship between ecology and modernism. He is a senior editor for *Methods in Ecology and Evolution*, and the author of several books, including *A Primer of Ecological Statistics*, *Stepping in the Same River Twice: Replication in Biological Research*, *Vanishing Point*, and *Carnivorous Plants: Physiology, Ecology, and Evolution*. Ellison received a PhD in ecology and evolutionary biology from Brown University. Contact him at aellison@fas.harvard.edu.

Margo Seltzer is the Herchel Smith Professor of Computer Science and the faculty director for the Center for Research on Computation and Society in Harvard's John A. Paulson School of Engineering and Applied Sciences. In September 2018, she will assume a Canada 150 Research Chair and the Cheriton Family Chair in Computer Systems at the University of British Columbia. Her research interests include systems, construed quite broadly: systems for capturing and accessing provenance, file systems, databases, transaction processing systems, storage and analysis of graph-structured data, new architectures for parallelizing execution, and systems that apply technology to problems in healthcare. Seltzer received a PhD in computer science from the University of California at Berkeley. Contact her at margo@eecs.harvard.edu.