

CPR - A Comprehensive Provenance Record for Verification Workflows in Whole Tale

No Author Given

No Institute Given

1 Introduction

An increasing number of journal publishers verify computational artifacts as part of the peer-review process. Although the problems of defining and achieving computational reproducibility have proven thorny generally, the particular problems publishers aim to detect in this context are well defined. Questions representative publishers answer via verification workflows include:

- Is the description in the text and supplementary materials sufficient to enable others to repeat the reported computations?
- Does repeating the computations yield the reported results?

Platforms such as Binder [2] and Whole Tale [1] provide environments for assessing reproducibility of computational artifacts by these standards via what is essentially *black-box testing* of the computational workflow. A verifier uses information provided in the paper to (1) set up the required computational environment; (2) stage input data; (3) trigger a sequence of automated computations; and (4) allow these computations to run to completion. Finally the verifier confirms that the products of the computations match the description in the paper.

Whole Tale further aims to enable verifiers to observe *how* automated computational workflows produce intermediate and final artifacts. Ultimately this will allow publishers to ask a third general question:

- Is the authors' description of the roles played by various software components consistent with the observed flow of data through those components?

This capability will provide verifiers means for *white-box* testing of the computations reported in a paper. Specifically, it will enable a verifier to detect cases where the sequence of computational steps and flow of data between these steps does not conform to the description given in the paper. The demonstration described here exercises and demonstrates capabilities of the tools Whole Tale employs to record, store, query, and visualize the flow of data through computational workflows for this purpose.

2 The CPR Toolkit

The CPR (Comprehensive Provenance Record) Toolkit is a suite of tools for recording, storing, querying, and visualizing the run-time provenance of artifacts produced by a run of a computational workflow. While the primary purpose of CPR is to automate the monitoring and management of provenance-relevant events and records associated with a Whole Tale *recorded run*, the toolkit can be employed in any Linux-based computing environment.

Figure 1 illustrates the flow of information through elements of the CPR toolkit. CPR employs `ReproZip` to observe system calls invoked as part of the recorded run and to record metadata about (1) the operating-system level processes comprising the overall computation; (2) the files accessed by these processes; and (3) the access mode for file accesses, i.e. whether the process opened the file for reading, writing, or both. `ReproZip` captures and records all of this information in a SQLite database with a schema specific to `ReproZip`.

Once a recorded run is complete, the `cpr` utility extracts these OS-level records from the `ReproZip` trace, transform them into RDF triples, and loads the triples into an RDF dataset in an instance of `Blazegraph`. The triples are expressed using a vocabulary developed to represent provenance information in the context of Whole Tale recorded run executions. The CPR vocabulary extends `PROV` and `ProvONE` with subclasses and new concepts specific to Whole Tale, supporting storage and query of provenance captured from multiple recorded runs and versions of multiple Tales (Figure 2). CPR can represent this vocabulary either as Datalog facts or as RDF triples.

The CPR vocabulary distinguishes between a number of general roles that files accessed during a run may play. A simple YAML file is used to declare prospectively the role of individual files, the contents of particular directories, or of entire directory trees. By using these declarations while converting a `ReproZip` trace to the CPR vocabulary, CPR is able to distinguish data files of scientific significance from shared libraries provided by the operating system or software dependencies.

Finally, the Geist reporting tool is used to pose SPARQL queries against the `Blazegraph` instance, to format the query results as reports, and to create visualizations of query results using `Graphviz`. Geist queries, reports, and visualizations may be parameterized; in Whole Tale we plan to create a predefined set of reports and visualizations following each recorded run.

3 Demonstration

The CPR demo is provided as a Git repository and associated Docker image that enable the examples to be run on any Linux, macOS, or Windows-based system with Git, Docker, and GNU Make installed. Each example uses the CPR toolkit

to record OS-level provenance information from a run of different computational workflow, to load a Blazegraph instance with the resulting CPR traces, and to produce a set of reports and visualizations via SPARQL queries.

A Makefile in the top directory of the demo repository provides targets for pulling the Docker image from Dockerhub (`pull-image`), building the Docker image locally (`build-image`), for running the examples (`run-examples`), and for deleting all of the reports, visualizations, and other artifacts generated for each example (`clean-examples`). Because the expected results are included in the repository, successful reproduction of the example products is demonstrated by issuing the commands `make clean-examples` and `make run-examples` and confirming that `git diff` reports no differences.

3.1 Queries and Visualizations

Queries and visualizations produced for each example included the following. A question phrased in English summarizes the intent of each.

What files are employed as inputs and produced as outputs of the run?

What are the input and output data files for each process in the run?

Which files input to a run are used to produce a particular output file?

Which output artifacts are affected by a particular input file?

What programs and script invocations occur as part of the run?

What programs and script invocations contribute to the production of a particular output artifact?

What constraints on the order of execution of different processes does the observed flow of data imply?

3.2 Observations

A key challenge in making provenance useful to researchers and verifiers alike is highlighting the small subset of recorded events of direct relevance to the scientific purpose of an execution.

An execution of a one-line Python 3 script that prints "Hello World" can involve reading as many as X different files from disk in addition to the users' single-line Python file.

CPR minimizes such provenance "noise" using a CPR profile that assigns distinct roles to files loaded from particular locations on the system.

References

1. Brinckman, A., Chard, K., Gaffney, N., Hategan, M., Jones, M.B., Kowalik, K., Kulasekaran, S., Ludäscher, B., Mecum, B.D., Nabrzyski, J., Stodden, V., Taylor, I.J., Turk, M.J., Turner, K.: Computing Environments for Reproducibility: Capturing the “Whole Tale”. *FGCS* **94**, 854–867 (2019). <https://doi.org/10.1016/j.future.2017.12.029>, <http://www.sciencedirect.com/science/article/pii/S0167739X17310695>
2. Jupyter-Project: Binder 2.0 - Reproducible, Interactive, Sharable Environments for Science at Scale. 17th Python in Science Conference (2018)