

Wanted: standards for automatic reproducibility of computational experiments

SAMUEL GRAYSON, University of Illinois Urbana-Champaign, USA

JOSHUA TEVES, Sandia National Laboratories, USA

REED MILEWICZ, Sandia National Laboratories, USA

DANIEL S. KATZ, University of Illinois Urbana-Champaign Department of Computer Science, USA

DARKO MARINOV, University of Illinois Urbana-Champaign, USA

1 INTRODUCTION

A computational experiment is reproducible if another team using the same experimental infrastructure can make a measurement that concurs with the original. In practice, reproducers will still need to look at the code by hand to see how to build necessary libraries, configure parameters, find data, and invoke the experiment; it is not *automatic*. To enable automatic reproducibility, one would need a description language which could list the relevant commands with machine-readable metadata attached to describe what the command does. It is not enough for the language to merely contain this command in a heap of other commands; e.g., a Makefile which defines a rule for executing the experiment alongside rules for compiling intermediate pieces is not sufficient, because there is no machine-readable way to know which of the Make rules executes the experiment.

Automatically identifying the “main” command which executes the experiment is critical for:

- **Artifact evaluators:** With manual reproducibility, artifact evaluators spend time learning how to set up, configure, build, and review the artifact. Automatic reproducibility would be the canonical place for experiment computational scientists to concisely communicate these steps. Unlike `README.txt` it would be human- and machine-readable.
- **Users seeking to re-execute with different parameters:** With manual reproducibility, users have to dig through the experiment’s documentation or, more likely, source code to discover how to supply parameters. Automatic reproducibility can also specify how to set these parameters.
- **Large-scale re-execution experiments:** Collberg and Proebsting [1] do a large-scale study of repeatability of computational experiments in computer science with manual effort. While their results are seminal, it is difficult to repeat in other domains or extend that experiment without spending a huge amount of human-hours figuring out how to run experiments. If Collberg and Proebsting or some other software-engineering researchers *do* embark to figure out how to run a certain experiment, there is standardized way for them to share their steps with other researchers.

2 HOW TO GET AUTOMATIC REPRODUCIBILITY

This is not the final proposal for the complete vocabulary; the peer-review process is not well-suited to iterate on technical details. The point of this article is to argue that the community should spend effort developing this vocabulary.

Unpublished working draft. Not for distribution.

2.1 Semantic web description

This language could be implemented as a vocabulary for linked data in the semantic web. Linked data is preferable for these reasons:

1. Linked data is open to extensions.
2. It is possible to link to other resources in linked data.
3. There is already a rich set of ontologies for describing digital and physical resources (RO-crate, wf4prov, software project description, scientific hypotheses, CiTO) in linked data.
4. There is already a rich ecosystem for authoring ontologies and validating documents within those ontologies.

Linked data is already used for other long-term preservation standards, such as RO-crate. The template of RDF/XML looks like this:

```
<?xml version="1.0"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:cito="http://purl.org/spar/cito"
  xmlns:doco="http://purl.org/spar/doco/2015-07-03"
  xmlns:prov="https://www.w3.org/TR/2013/PR-prov-o-20130312/"
  xmlns:wfdesc="http://purl.org/wf4ever/wfdesc#"
  xmlns="http://example.org/execution-description/1.0" >
...
</rdf:RDF>
```

According to the RDF/XML specification, This imports several other vocabularies behind a namespace. E.g., `rdf:` type refers to type in the `rdf` namespace, which points to `http://www.w3.org/1999/02/22-rdf-syntax-ns#`. XML tags with no namespace are resolved within the default namespace, which is our proposed execution-description vocabulary.

2.2 Language description

At a very basic level, one could have commands and the purpose that they serve:

```
<process>
  <command>./execute --input data.csv</command>
  <purpose>generates data</purpose>
</process>
<process>
  <command>python3 main.py</command>
  <purpose>plots figures</purpose>
</process>
```

While we could define conventions around what to name the content of the “purpose” tag, it would be more powerful if the language could link directly to the claims in the publication. The CiTO vocabulary [5] already defines a vocabulary for describing citations.

```
<process>
  <command>./execute --input data.csv</command>
  <purpose>
    <!-- links to an entire publication -->
```

```

105     <cito:isCitedAsEvidenceBy rdf:resource="https://doi.org/10.1234/123456789" />
106   </purpose>
107 </process>
108

```

If the publisher hosts an RDF description at the URL “https://doi.org/10.1234/123456789” when the HTTP request content-type header is application/rdf+xml, then this creates a web of linked data. The publisher may have the title, authors, date published, and other metadata using Dublin Core metadata terms, for example. This is the dream of linked data: machine-readable data by different authors hosted in different locations linking together seamlessly. Even if the publisher does not have an RDF+XML description, third parties can make claims about “https://doi.org/10.1234/123456789”, although those claims would not be as easily discoverable. One could even reference a Nanopublication, which is a semantic web description of the scientific claim.

The purpose description can be even more granular, using the DoCO vocabulary [2], which describes documents.

```

119 <purpose>
120   <prov:generated>
121     <doco:figure>
122       <dc:title>Figure 2b</dc:title>
123       <dc:isPartOf rdf:resource="https://doi.org/10.1234/123456789" />
124     </doco:figure>
125   </prov:generated>
126 </purpose>
127
128

```

With this complete, anyone should be able to execute the experiments which generate figures or claims in the paper if they are labeled by the computational scientist in this language.

One can view specifying the software environment as just prerequisite steps in the computational DAG.

```

133 <process>
134   <wfdesc:hasOutput>
135     <wfdesc:Output rdf:nodeID="conda-env-out" />
136   </wfdesc:hasOutput>
137   <command>conda env create --name experiment-123 --file environment.yml</command>
138 </process>
139 <process>
140   <wfdesc:hasInput>
141     <wfdesc:Input rdf:nodeID="conda-env-in" />
142   </wfdesc:hasInput>
143   <purpose>figure 4</purpose>
144   <command>conda run --name experiment-123 ./plot-figure.py</command>
145 </process>
146 <wfdesc:DataLink>
147   <wfdesc:hasSource rdf:nodeID="xy-dataset-out" />
148   <wfdesc:hasSink rdf:nodeID="xy-dataset-in" />
149 </wfdesc:DataLink>
150
151
152

```

Now a machine knows that the `conda env create ...` must be run before, and the script should be run within the resulting conda environment.

The purpose of an execution language is not to usurp the build-system or workflow engine, which both already handle task DAGs; there must be some minimal support for DAGs just for the cases where the DAG of tasks is not already encoded in a build-system or workflow engine.

In addition to specifying the computational environment and command to run, this language is an ideal candidate to also describe the parameters of the experiment, like: One could specify range of valid values or list options. With this complete, one can even do automated parameter-space search studies, multi-fidelity uncertainty quantification, automated outcome-preserving input minimization, and other automatic experiments.

Retrospective provenance seeks to encode how we got to a specific result. Developers can put summary statistics of intermediate results into the provenance description; if re-executions diverge, users can locate which stage amplifies error the most. A tool might use system-call interposition to learn about a processes reads, writes, and forks. This would be better at identifying and recording intermediate results. When reproducing some computational experiment, users can check the intermediate results to see where they begin to differ. Wfprov is one vocabulary for specifying retrospective provenance, and there is already an experimental plugin for Nextflow which targets wfprov [3].

Users may also want to know how much computational resources (CPU time, disk space, and RAM) the computational experiment requires. Provenance is the ideal place for computational experiments who already ran the experiment to put this information. This way, users seeking to reproduce the computational experiment know how many resources to request (ahead-of-time allocations are usually required for batch-scheduled machines).

3 MAKING EASY ON-RAMPS FOR ADOPTION

The execution language should seek to describe existing software frameworks, not replace them. In particular, execution should not replace workflow engines. They should instead be wrapped as process-nodes within the execution language. Computational experiments can continue using their existing build-system and workflow.

A execution description could even be “captured” from an interactive shell session with the user. They would invoke a shell that records every command, its exit status, its read-files, and its write-files (using syscall interposition). The user would execute their build-system like normal within this shell. When the user exits, the shell will create a DAG based on the read-files and write-files. For each output file that is not consumed by another command, the shell would prompt the user to describe that command’s purpose. Finally, the shell would output a execution description.

In cases where the description cannot be captured by an interactive shell session, one can encode the system of logic that humans would use to deduce the experimental structure. This is similar to the approach taken by FlaPy [4], a large-scale re-execution study for Python unittests, to install the Python environment. For example,

1. If `shell.nix`, `flake.nix`, or `environment.yml` exists, then use Nix, Nix flakes, or Conda as the environment for future commands.
2. If the environment is not already set, and `requirements.txt` exists, use Pip and Virtualenv as the environment for future commands.
3. If `CMakeLists.txt` exists, set `cmake` and `make` as commands.
4. If `configure.sh` exists, set `./configure` and `make` as command.
5. If `Makefile` exists, set `make` as command. ...

In the best case, the execution description would be uploaded to the same repository or location that contains the source for the computational experiment. This way it can be maintained and used by the original developers, and it is easily discoverable by users. Alternatively, execution descriptions could be placed in a central execution-description

repository owned by software-engineering researchers. The execution description is linked data, so it can live anywhere, and existing strategies for finding, filtering, and trusting linked data sets would work for execution descriptions.

This is similar to the approach taken by Python for type annotations. Type annotations are easiest to maintain in the original repository, but if the original repository rejects type annotations, there is still a home for them in typeshed.

4 IS THIS ANOTHER COMPETING STANDARD?

It is something extra users have to do, but it does not attempt to displace their existing practice.

5 INCENTIVES FOR COMPUTATIONAL SCIENTISTS TO MAINTAIN EXECUTION DESCRIPTIONS

Computational scientists are incentivized to describe their project this way to benefit from the work of software-engineering researchers. When software-engineering researchers do a large-scale execution study, it is a “free” reproduction of their work.

To get an artifact evaluation badge, normally computational scientists would have to write a natural language description of what the software environment, what the commands are, how to run them, and where does the data end up. The artifact evaluator has to read, interpret, and execute their description by hand. An execution description could make this nearly automatic; if a execution description exists, the artifact evaluator uses an executor which understands the language and runs all of the commands that reference the manuscript in their purpose tag. The only manual labor is comparing these results to those in the paper. Even that comparison can be simplified, if the last step in the execution description outputs a boolean representing “is the hypothesis proven?”; the reviewer just needs to see that all of these output “true”.

REFERENCES

- [1] Christian Collberg and Todd A. Proebsting. 2016. Repeatability in computer systems research. en. *Communications of the ACM*, 59, 3, (Feb. 2016), 62–69. doi: 10.1145/2812803.
- [2] Alexandru Constantin, Silvio Peroni, Steve Pettifer, David Shotton, and Fabio Vitali. 2016. The Document Components Ontology (DoCO). en. *Semantic Web*, 7, 2, (Jan. 2016), 167–181. Publisher: IOS Press. doi: 10.3233/SW-150177.
- [3] Bruno Grande, Ben Sherman, and Paolo Di Tomasso. 2023. Nf-prov. original-date: 2022-12-19T21:16:30Z. (May 2023). Retrieved May 25, 2023 from <https://github.com/Sage-Bionetworks-Workflows/nf-prov>.
- [4] Martin Gruber, Stephan Lukasczyk, Florian Kroiß, and Gordon Fraser. 2021. An Empirical Study of Flaky Tests in Python. In *2021 14th IEEE Conference on Software Testing, Verification and Validation (ICST)*, ISSN: 2159-4848. (Apr. 2021), 148–158. doi: 10.1109/ICST49551.2021.00026.
- [5] David Shotton. 2010. CiTO, the Citation Typing Ontology. en. *Journal of Biomedical Semantics*, 1, 1, (June 2010), S6. doi: 10.1186/2041-1480-1-S1-S6.