

Reviewer Report

Title: Sharing interoperable workflow provenance: A review of best practices and their practical application in CWLProv

Version: Original Submission **Date: 1/22/2019**

Reviewer name: Alban Gaignard

Reviewer Comments to Author:

Review "Sharing interoperable workflow provenance: a review of best practices and their practical application in CWLProv"

This paper proposes an in-depth, extensive, literature review in the field of reproducible computational sciences as well as a prototype implementation leveraging state-of-the-art standardization efforts. Main issues identified are the incompleteness of captured provenance information and the lack of interoperability. Even if CWL is emerging as a standard and system-independent language to represent scientific workflows, there is no standards to represent, share, and reproduce computational analysis as workflow executions. As a result of the literature review, the authors propose (1) a broad list of requirements for better reproducibility, and (2) a hierarchy of provenance levels addressing various needs, from workflow debug, workflow re-execution and packaging, towards better re-use and re-purposing for scientists, based on domain-specific annotated provenance. Then, the CWLProv implementation is described and evaluated on 3 real-life bioinformatics workflows.

My main concern regarding this work is that it is often stated that the re-usability of workflow resources (methods / input or output data) is facilitated but it is difficult to evaluate this claim based on CWLProv features and the proposed experiments. It is clear that re-execution of workflows is facilitated but it is unclear to what extent produced/analysed data can be considered for secondary use. In addition, the "pragmatic" interoperability should refer to top-level provenance and thus domain-specific annotations referring to the scientific context of the computational experiment. The experiments don't clearly show how CWLprov goes into the direction of (still ambitious and challenging) domain-specific provenance.

I've also a technical concern regarding the FAIRness of the approach since some of the requirements could be addressed following the (5-star) Linked Data principles. This point should be addressed in the discussion.

Finally, I tried to browse the research objects provided as supporting material but unfortunately I could not access the resource. Logs are provided at the end of the review.

Apart from that, the paper is well written (sometime a bit long) and well illustrated. The structure is easy to follow. This paper results from an impressive and high-quality work for a very timely topic, with urgent community needs. It should be published, after minor clarifications and corrections.

Detailed comments

Introduction

In Key Points, 4th point, a space is missing in "CWLProvoutcome"

Background and related works

Methods

Are the methods appropriate to the aims of the study, are they well described, and are necessary controls included? Choose an item.

Conclusions

Are the conclusions adequately supported by the data shown? Choose an item.

Reporting Standards

Does the manuscript adhere to the journal's guidelines on [minimum standards of reporting?](#) Choose an item.

Choose an item.

Statistics

Are you able to assess all statistics in the manuscript, including the appropriateness of statistical tests used? Choose an item.

Quality of Written English

Please indicate the quality of language in the manuscript: Choose an item.

Declaration of Competing Interests

Please complete a declaration of competing interests, considering the following questions:

- Have you in the past five years received reimbursements, fees, funding, or salary from an organisation that may in any way gain or lose financially from the publication of this manuscript, either now or in the future?
- Do you hold any stocks or shares in an organisation that may in any way gain or lose financially from the publication of this manuscript, either now or in the future?
- Do you hold or are you currently applying for any patents relating to the content of the manuscript?
- Have you received reimbursements, fees, funding, or salary from an organization that holds or has applied for patents relating to the content of the manuscript?
- Do you have any other financial competing interests?
- Do you have any non-financial competing interests in relation to this paper?

If you can answer no to all of the above, write 'I declare that I have no competing interests' below. If your reply is yes to any, please give details below.

I declare that I have no competing interests.

I agree to the open peer review policy of the journal. I understand that my name will be included on my report to the authors and, if the manuscript is accepted for publication, my named report including any attachments I upload will be posted on the website along with the authors' responses. I agree for my report to be made available under an Open Access Creative Commons CC-BY license (<http://creativecommons.org/licenses/by/4.0/>). I understand that any comments which I do not wish to be included in my named report can be included as confidential comments to the editors, which will not be published.

Choose an item.

To further support our reviewers, we have joined with Publons, where you can gain additional credit to further highlight your hard work (see: <https://publons.com/journal/530/gigascience>). On publication of this paper, your review will be automatically added to Publons, you can then choose whether or not to claim your Publons credit. I understand this statement.

Yes Choose an item.