

A PROV-Compliant Approach for the Script-to-Workflow Process

Lucas A. M. C. Carvalho ^{a,*},
Khalid Belhajjame ^b
and Claudia Bauzer Medeiros ^a

^a *Institute of Computing, University of Campinas, Campinas, Brazil*

E-mails: lucas.carvalho@ic.unicamp.br, cmbm@ic.unicamp.br

^b *LAMSADE, Paris-Dauphine University, Paris, France*

E-mail: khalid.belhajjame@dauphine.fr

Abstract. Scientific discovery and analysis are increasingly computational and data-driven. Scripting languages, such as Shell, Python and R, are the means of choice of the majority of scientists to encode and run their simulations and data analyses. Although widely used, scripts are hard to understand, adapt, reuse, and reproduce. To tackle the problems faced by scripts, several approaches have been proposed such as YesWorkflow and noWorkflow. However, they neither allow to fully document the experiment nor do they help when third parties want to reuse just part of the code. Scientific Workflow Management Systems (SWfMSs) are being increasingly recognized to mitigate these problems. They help to document and reuse experiments by supporting scientists in the design and execution of their experiments, which are specified and run as interconnected (reusable) workflow components (a.k.a. building blocks). Taking this into account, we designed W2Share, a novel approach for the management, reuse, and reproducibility of script-based experiments. W2Share transforms a script into an executable workflow that is accompanied by annotations, example datasets and provenance traces of their execution, all of which encapsulated into a workflow research object. This allows third party users to understand the data analysis encoded by the original script, run the associated workflow using the same or different datasets, or even repurpose it for a different analysis. W2Share also enables traceability of the script-to-workflow process, thereby establishing trust in this process. All processes in W2Share follow a methodology that is based on requirements that we elicited for this purpose. The methodology exploits tools and standards that have been developed by the scientific community, in particular, YesWorkflow, Research Objects and the W3C PROV. This paper highlights the main components of W2Share, which is showcased through a real world use case from Molecular Dynamics. We furthermore validate our approach by testing the ability to answer competency questions that address the script-to-workflow process.

Keywords: Scientific Workflows, Script-to-Reproducible Research, Workflow Research Objects, Provenance

1. Introduction

Scripts and Scientific Workflow Management Systems (SWfMSs) [1, 2] are common approaches that have been used to automate the execution flow of processes and data analysis in scientific (computational) experiments¹. Scripts are widely adopted in many disciplines to create pipelines for experiment ex-

ecution, e.g., to clean and analyze a large amount of data. However, they are hard to understand, adapt, and reuse, often containing hundreds of lines of domain-specific code. This, in turn, forces scientists to repeatedly (re)code scripts that perform the same functions, since the effort to reuse is not worthwhile, and reproducibility is restricted to repeating the execution of exactly the same script. For this reason, several solutions have been proposed to aid experiment reproducibility for script-based environments such as Jupyter Notebooks [3], ReproZip [4], YesWorkflow [5], and noWorkflow [6].

*Lucas A. M. C. Carvalho. E-mail: lucas.carvalho@ic.unicamp.br

¹In this paper, the term *experiment* refers to scientific experiments that are executed *in silico* – e.g., simulations.

1 Though those solutions help scientists capture exper-
2 imental details, they neither allow to fully docu-
3 ment the experiment, nor do they help when third par-
4 ties want to reuse just part of the code. For exam-
5 ple, the workflow-like graph obtained using YesWork-
6 flow is abstract (in the sense that it cannot be executed
7 by the scientists). On the other hand, the provenance
8 traces captured using noWorkflow are fine-grained,
9 and therefore cumbersome for the user who would
10 like to understand the lineage of the script results [7].
11 SWfMSs [8], on the other hand, help documentation
12 and reuse by supporting scientists in the design and ex-
13 ecution of their experiments, which are specified and
14 run as interconnected (reusable) workflow components
15 (a.k.a. building blocks).

16 While workflows are better than scripts for under-
17 standability and reuse, they still require additional docu-
18 mentation to support reproducibility. To this end, we
19 designed and implemented W2Share, a computational
20 framework that supports a (script-to-reproducible re-
21 search) methodology. The methodology, implemented
22 in W2Share via a suite of tools, guides scientists
23 in a principled manner to transform scripts into re-
24 producible and reusable workflow research objects
25 (WRO) [9]; it drives the development of research ob-
26 jects that contain the scripts that the scientist authored
27 together with executable workflows that embody and
28 refine the computational analyses carried out by these
29 scripts and all associated data and documentations.
30 Our methodology thus leverages the concept of Work-
31 flow Research Objects as a means to ensure repro-
32 ducibility.

33 W2Share’s WRO encompasses information such as
34 the workflow itself, and datasets and provenance traces
35 related to its execution. The WRO model [9] allows the
36 aggregation of resources, explicitly specifying the re-
37 lationship between these resources and a workflow, us-
38 ing a suite of ontologies. W2Share’s WROs allow sci-
39 entists to understand the relationships between an ini-
40 tial script and the resulting workflow, and to document
41 workflows runs – e.g., annotations to describe the oper-
42 ations performed by the workflow, or links to other re-
43 sources, such as the provenance of the results obtained
44 by executing the workflow. Using W2Share, scientists
45 can share and reuse scripts through the corresponding
46 WROs.

47 Our approach differs in several ways from similar
48 work to convert scripts into workflows such as [10–
49 12]. In particular, our steps to convert scripts into ex-
50 ecutable workflows are more generic, in the sense that
51 they are independent from the script language and

1 the workflow system, while these other solutions are
2 mainly designed for specific environments. Moreover,
3 we are not only concerned about the workflow specifi-
4 cation derived from the script code, but also to preserve
5 the script in the WRO, allowing scientists to check ex-
6 periment provenance and reproducibility.

7 For that, we support two forms of provenance [13]:
8 (1) prospective and (2) retrospective. Prospective
9 provenance captures the specification of steps and
10 their data dependencies for a given computational task
11 (whether it is a script or a workflow). Retrospective
12 provenance captures the steps executed and the order
13 of execution, along with the data consumed and pro-
14 duced by each step as well as different kind of meta-
15 data that help understanding and reproducing the exe-
16 cution.

17 The main contributions of this paper therefore in-
18 clude:

- 19 1. A methodology to guide scientists in a principled
20 manner to transform scripts into reproducible
21 and reusable workflow research objects²;
- 22 2. A data model that identifies the main elements of
23 the methodology and their relationships, which
24 helps automate the steps of the methodology, and
25 their documentation;
- 26 3. A computational framework that provides sci-
27 entists with the tooling necessary for (semi)-
28 automatically performing some of the steps of
29 the methodology. Here, we emphasize that we
30 make use of semantic web technologies, web
31 standards and tools developed by the scientific
32 community;
- 33 4. A case study to showcase our solution; and
- 34 5. An evaluation of the proposed model and conver-
35 sion mechanism via the identification and execu-
36 tion of competency questions.

37 This paper extends previous work [14] of ours,
38 where we defined and presented the methodology, and
39 showcased its use via a case study (through man-
40 ual implementation of the conversion). This paper de-
41 scribes how we now enable (semi) automatic script-
42 to-workflow conversion, and its validation using com-
43 petency queries that address requirements used to de-
44 sign our methodology. The automation of the conver-
45 sion process is based on our data model that iden-
46 tifies the resources that are used and generated by
47 our methodology as well as the agents responsible
48 for the execution.

49 ²This methodology was described in [14], and is included here for
50 completeness.
51

1 for performing the methodology steps. Last but not
2 least, besides providing traceability for experiment ex-
3 ecution (via workflow mechanisms), we innovate by
4 providing traceability for the script-to-workflow pro-
5 cess itself. We describe the design and implementation
6 of this extended conversion process, which takes ad-
7 vantage of ontologies adopted by the scientific com-
8 munity, namely W3C PROV-O [15], Web Annotation
9 Data Model³, and Research Object ontology. These
10 ontologies support the semantics of the traceability of
11 the *script-to-workflow* process.

12 This paper is structured as follows. Section 2 in-
13 troduces our methodology specification. Section 3
14 presents the model that describes its main elements.
15 Section 4 describes the case study and the implemen-
16 tation of the methodology steps. Our conversion steps
17 are described in Sections 5 and 6. Section 5 describes
18 the first step, showing how we map scripts to abstract
19 workflows using ontologies. Section 6 shows how an
20 executable workflow is generated from abstract work-
21 flows. Section 7 presents the evaluation of our pro-
22 posed approach. Section 8 discusses related work. Sec-
23 tion 9 summarizes our results and identifies future
24 work.

26 2. Methodology for Script Conversion into WRO

27
28
29 Parts of sections 2.1 and 4 appeared in our paper
30 [14] in which we defined our methodology and exem-
31 plified its application. This has been included in this
32 paper for clarity sake, and to make it self-contained.

33 2.1. Overview

34
35
36 Our methodology guides scientists throughout the
37 conversion of script-based experiments into executable
38 workflows, and then packaging all the resources and
39 annotations into a Workflow Research Object (WRO)
40 [16]. WROs are the means through which experiments
41 can be reused, audited and documented. As mentioned
42 in Section 1, our WRO encapsulates the scripts and
43 the corresponding executable workflows together with
44 other resources, such as datasets and provenance traces
45 of their execution.

46 The methodology was designed to meet five ma-
47 jor requirements that were derived during our long-
48 time collaboration with scientists that run scripts for
49 their computational experiments (e.g., in bioinformat-

1 ics, chemistry and agriculture). These requirements are
2 the following:

3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
Requirement 1 Let S be a script that embodies a computational experiment. The scientist needs a view of S that identifies the main processing units and dependencies between such processing units.

This helps the scientist understand S , and the main processing units that are relevant from the point of view of the *in silico* analysis implemented by the script, as well as the dependencies between such units.

We call this view an *abstract workflow*. In more detail, an abstract workflow, for the purposes of this paper, is a process, in which the steps designate script blocks, and the dependencies designate data dependencies between these blocks. The workflow is abstract in that it is not executable *per se*, but rather provides a process view of a script at a higher level of granularity (logical steps as opposed to script instructions).

Requirement 2 The scientist should be able to execute the workflow that embodies the script S .

Though seemingly obvious, this is far from being a trivial requirement. It is not enough to "be able to execute". This execution should reflect what is done in script S . In other words, not only should the workflow generated be executable; the scientist must be given the means to compare its results to those of script execution, and validating the workflow as a valid embodiment of the script.

Requirement 3 The scientist should be able to modify the workflow that embodies script S , to use different computational and data resources.

Not only may a scientist be able to replicate the computational experiment encoded by S ; s/he may want to repeat the analysis implemented in the script using third party resources.

The new (modified) workflow(s) correspond to variants of the initial workflow. They will help the user, for example, to inspect if the results obtained by script S can be reproduced using different resources (algorithms and datasets). Scientists will also be able to compare the execution of S with that of the variants (e.g., if web services are invoked instead of a local code implementation).

Requirement 4 Provenance information should be recorded.

Provenance information is key to traceability and quality assessment. This involves not only the provenance obtained by workflow execution. This requirement also implies recording the transformations car-

51 ³<https://w3.org/TR/annotation-model>

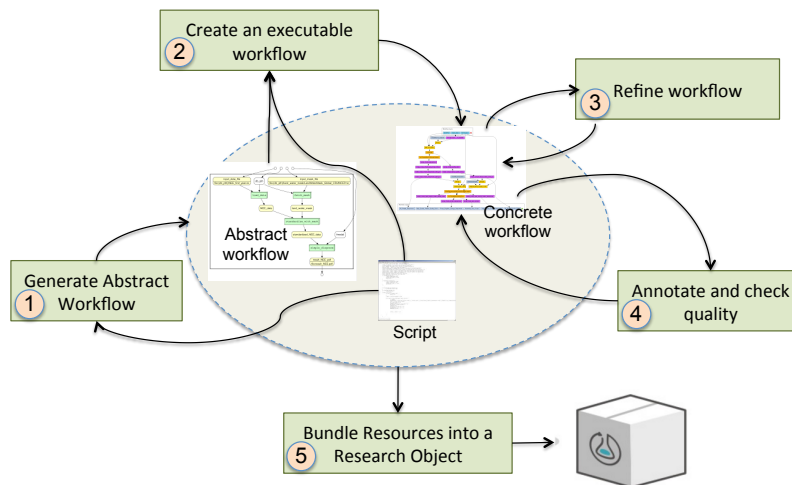


Fig. 1. Methodology for converting scripts into reproducible Workflow Research Objects, extracted from [14].

ried out to transform the script into a workflow that embodies the script. Moreover, the workflow variants also need to be recorded. As stressed by [17], provenance that is provided by the execution of a workflow corresponds to a workflow trace, and can be used for several purposes, such as to support dynamic steering of workflows [18, 19].

Requirement 5 All elements necessary to reproduce the experiment need to be captured together to promote reproducibility.

We follow the definition of [17]: "reproducibility denotes the ability for a third party who has access to the description of the original experiment and its results to reproduce those results using a possibly different setting, with the goal of confirming or disputing the original experimenter's claims." Missier et al. [17] also differentiate reproducibility from repeatability, for which results must be the same, and no changes are made anywhere.

Full reproducibility and reusability require ensuring that all elements of an experiment are recorded. S , the initial workflow, and all of its variants should be made available together with auxiliary resources that will allow understanding how these workflows came to be, and where they should be used.

Given the five requirements, we proposed our methodology composed by five inter-related steps (see figure 1). Step 1, **Generate abstract workflow**, is used to produce an abstract workflow W_a based on a script S provided by a user. This stage elicits the main processing units that constitute the analysis implemented by the script, and their data dependencies. This requires

user intervention, to identify these units and dependencies within the script. Workflow W_a obtained from Step 1 is abstract in the sense that it cannot be executed – it is only a workflow-like high level specification of the script. Already at this stage, even though unable to execute the workflow, this is already a step towards promoting understandability – the abstract workflow is a high-level specification of the script, and can be visualized as a graph linking computational units. The objective of this phase is to address Requirement 1.

Figure 2 illustrates this step. The left side shows an excerpt of the script (from hundreds of lines of script code) and the right side the corresponding abstract workflow. This example will be discussed at length in subsequent sections; the figure is introduced here to give a high level view of this first step of the conversion process.

Step 2, **Create an executable workflow**, converts the abstract workflow W_a into an executable one W_e . The objective of this phase is to address Requirement 2. This is achieved by actually replacing each processing unit in the abstract workflow by its implementation (e.g., encapsulating the corresponding script code), and adding code to allow the required I/O operations across these units.

Scientists frequently try different variants of a computational experiment, e.g., to improve results, or to check alternatives. Script-based experiments are not easily modified, and it is hard to keep track of these variants. Tools such as version control systems allow to track the versions/changes of scripts and programs in general. Our methodology contemplates this activity. Step 3, **Refine workflow**, addresses Requirement 3

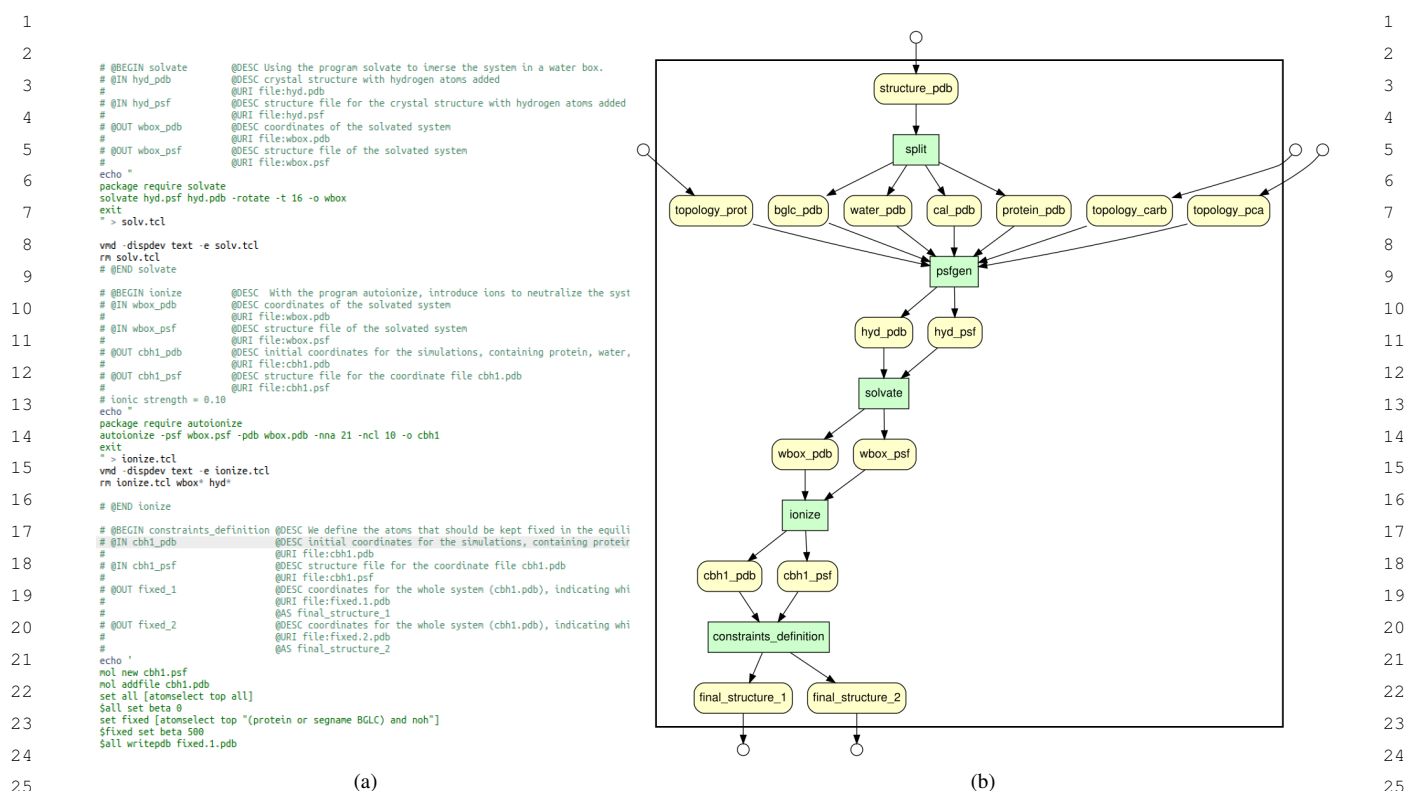


Fig. 2. The first step of our methodology concerns the generation of an abstract workflow from the script, illustrated here by (a) excerpt of the script and (b) the corresponding abstract workflow. This enhances understandability, and thus collaborative work.

and supports full reusability. It allows the creation of variants of the executable workflow, e.g., by adding new processing units, or changing data sources.

At the end of Step 3, the scientist will have one or more workflow variants $We_1 \dots We_n$. The idea, here, is that there is a difference between the concepts of repeatability and reproducibility. The first consists in exact reproduction of the experiment – running the same code, with the same data sets. Reproducibility, on the other hand, means that the results of an experiment should be reproducible, but not necessarily by invoking the same processes – e.g., code optimization can improve execution time.

Moreover, versioning allows scientists to try out variants of an experiment [20, 21], comparing and testing alternative outcomes. Overall, there are several kinds of refinements that can be performed at this step, all of which facilitated by the use of workflows and their components as reusable units.

During steps 2 and 3, provenance data both from the workflow executions and the process of conversion are collected to be used in Steps 4 and 5, and address requirement 4.

Step 4, **Annotate and check quality**, is in charge of evaluating whether the workflow reproduces the script results within some scientist-defined tolerance thresholds. It takes advantage of workflow execution mechanisms, that keep track of execution traces. Step 4 uses the workflow information generated in the previous steps, including provenance traces.

Finally, in Step 5, **Bundle Resources into a Research Object**, the workflow and the auxiliary resources, i.e., annotations, provenance traces, datasets, among others, are packaged into a WRO. WROs are then stored and made available to third parties for experiment validation, reproducibility checks, and reuse of workflow components. The objective of this phase is to address Requirement 5.

3. W2Share's Data Model: Supporting the Methodology

The full implementation of the methodology requires an appropriate data model, described here. It helps dynamic documentation of the conversion pro-

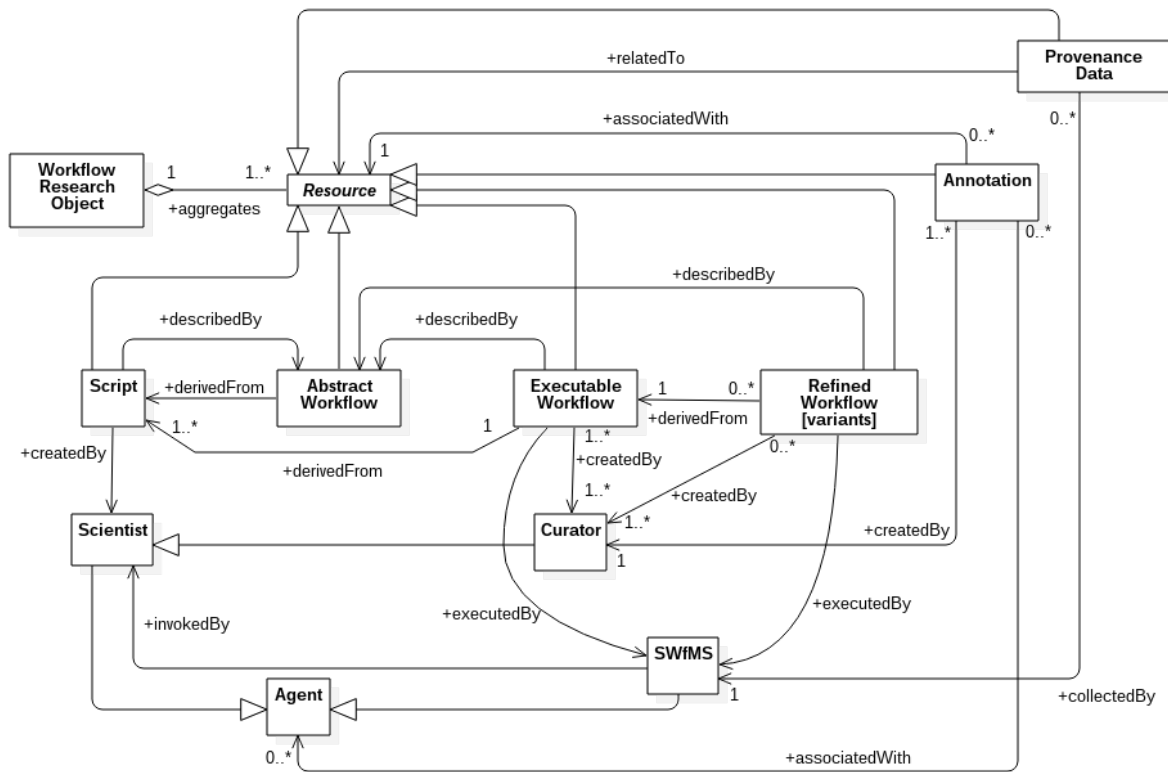


Fig. 3. Model describing main elements and relationships used in our methodology.

cess, thereby ensuring traceability of that process. We point out this adds a new dimension to traceability, which is usually restricted to the execution of experiments, but not to their evolution.

Figure 3 shows the UML class diagram of W2Share's data model, which reproduces the main entities and relationships involved in our methodology. The model describes the Resources that compose a WRO and the Agents that are responsible for creating these Resources.

Resources include, for instance, *Script*, *Abstract Workflow*, *Annotation*, *Provenance Data*, among others. These resources are not independent of one another – the model accounts for the relationships created in the transformation process – from *Script* to *Abstract Workflow* to *Executable workflow*, which can then give origin to several *Refined workflows* (variants).

Agents perform the activities in the methodology. As example of an Agent, a SWfMS (Scientific workflow Management System), which is invoked by a Curator (another Agent), is responsible for executing

workflows and collecting Provenance Data. The model differentiates between generic Scientists and Curators, scientists who are knowledgeable about documentation and resource management.

As mentioned before, we support two kinds of traceability – of experiment execution (based on SWfMS "logs") and of the conversion process itself. Traceability of the conversion process is enabled via relationships that are based on PROV. Examples include *annotatedBy*, *generatedBy*, *createdBy*, *derivedFrom*, and *collectedBy*. The adoption of PROV allows to navigate the derivation between the *Executable Workflow* and its variants in *Refined Workflow*. *Executable workflows* We are not derived from a Wa but directly from S . On the other hand, *Refined Workflows* $We_1, We_2 \dots We_n$ are derived from We .

Figure 4 shows a UML class diagram that refines part of figure 3, and is used to record the provenance of resources generated under our methodology. Here, we can see that *Scripts* are composed of *Code Blocks*; workflows (*Abstract Workflow*, *Executable Workflow* and *Refined Workflow* classes) are

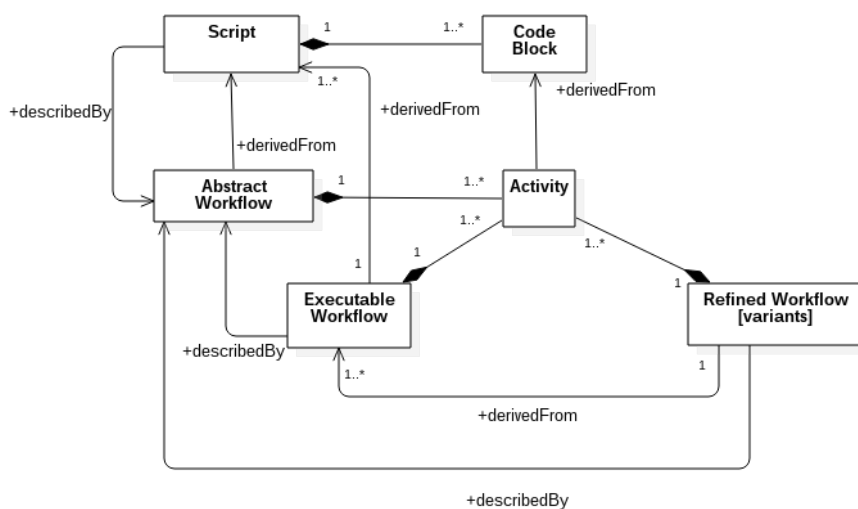


Fig. 4. Relationships between main entities regarding the tracking of the script-to-workflow conversion in our model.

composed of Activities, and the latter may be derived from script Code Blocks. Abstract Workflows describe Scripts and Workflows, thus allowing scientists to pose queries to explore prospective provenance. Eventual one-to-one relationship cardinalities were omitted from both figures, for readability.

Next, we show how to take advantage of the methodology to convert a script into a WRO in our case study from Molecular Dynamics.

4. Case Study – Molecular Dynamics

4.1. Overview

Molecular dynamics (MD) simulations consist of a series of algorithms developed to iteratively solve the set of coupled differential equations that determine the trajectories of individual atoms that constitute the particular physical system. This involves a long sequence of scripts and codes.

MD simulations are used in many branches of material sciences, computational engineering, physics and chemistry. A typical MD simulation experiment receives as input the structure, topology and force fields of the molecular system and produces molecular trajectories as output. Simulations are subject to a suite of parameters, including thermodynamic variables.

Many groups have implemented their specific MD simulations using special purpose scripts. In our case study, a suite of scripts was designed by physio-

chemists [22]; its inputs are the protein structure (obtained from the RCSB PDB protein data bank⁴), the simulation parameters and force field files.

There are many kinds of input files and variables, and their configuration varies with simulation processes. For instance, the input multimolecular structure contains the initial set of Cartesian coordinates for every atom/particle in the system, which will evolve in time in the MD simulation. This initial structure varies according to the system to be simulated and research area. Our case study requires immersing proteins in a solvent. Protein Cartesian atomic coordinates are made available in specialized data repositories, most notably the Protein Data Bank (PDB). Typical systems contain from several thousands to millions of covalently bound atoms.

Parts of the text in this section are based on [14].

4.2. Implementation of Methodology Steps

W2Share was conceived to take advantage of tools and standards that have been developed by the scientific community to support reproducibility and reuse, in particular YesWorkflow [5] and Research Objects. Its implementation includes elements of the PROV ontology, thereby facilitating provenance annotation. This is presented in Section 5 of this paper, and is one of the paper's contributions.

⁴<http://www.rcsb.org/pdb/>

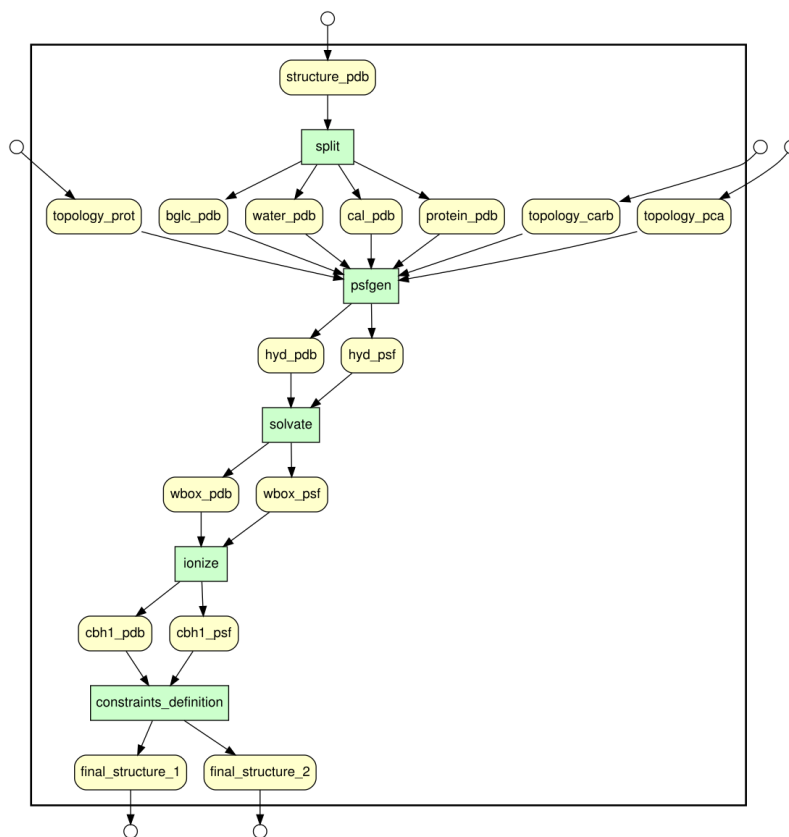


Fig. 5. Visualization of the abstract workflow of our MD case study, extracted from [14].

Consider the script that sets up an MD simulation. First, a scientist identifies the main processing units, and their dependencies. To do so, W2Share adopts the YesWorkflow tool. It enables scientists to annotate existing scripts with special comments that reveal the computational modules and data flows implicit in these scripts. YesWorkflow extracts and analyzes these comments, represents the scripts in terms of entities based on the typical scientific workflow model, and provides graphical renderings of this workflow. It does so by processing scientist-provided tags of the form `@tag value`, where `@tag` is a keyword that is recognized by YesWorkflow, and `value` is a value assigned to the tag. Tag recognition is script-language independent, therefore allowing a wide range of script-based experiments to be converted into a workflow representation and consequently a wider adoption of our methodology. W2Share creates the corresponding abstract workflow Wa (see figure 5 for the corresponding visualization) from the annotated script S , available in Listing A.1 in Appendix A.

This abstract workflow is a first approximation of what is needed for full reproducibility. Section 5 will detail how W2Share supports the creation of PROV-compliant machine-readable abstract workflows. The rest of this section will ignore these details, since they are not necessary to describe the full implementation of the methodology steps.

Given the abstract workflow Wa generated previously, the scientist needs to create an executable workflow We that embodies the data analysis and processes as depicted by Wa – and thus embodies the original script. For this, the scientist needs to specify, for each activity in the abstract workflow, the corresponding concrete activity that implements it. A simple, yet effective approach to do so consists in exploiting a readily available resource, namely the script code itself. Given an activity in Wa , the corresponding code in We can be generated by reusing the chunk (block) of the script that is associated with that abstract workflow activity. This approach for conversion comes with two advantages: (i) ease of conversion, since we are using a readily available resource, i.e. the script code,

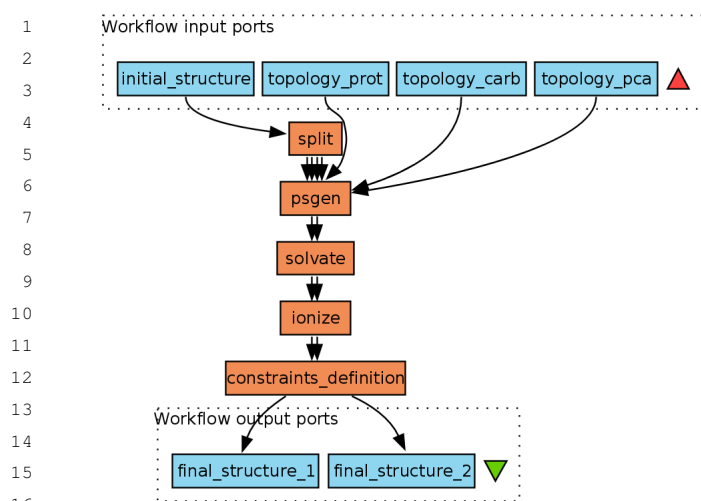


Fig. 6. Executable workflow of our MD case study, extracted from [14].

and (ii) the ability to check and debug the execution of *We* against the script execution, to correct eventual mistakes in script-to-workflow conversion. Once the scientist specifies the implementation of each activity in *Wa*, a concrete workflow specification *We* that is conform to a given scientific workflow system can be created. This manual conversion, detailed at length in [14], is now supported by W2Share – the semi-automatic implementation of Step 2 is detailed in Section 6.

Without loss of generality, we used the Taverna system [23], although our solution can be adapted to other scientific workflow systems. We chose Taverna as our implementation platform due to its widespread adoption in several eScience domains and because it supports the execution of shell scripts, the script language adopted in our case study. Figure 6 shows the executable workflow *We*, which is derived from *S* and described by *Wa*.

W2Share also helps scientists in creating workflow variants. For instance, in our case study, scripts use local data files containing protein coordinates which scientists download from authoritative web sources. This forces them to download such files from the web, and to update them locally whenever they are modified, moreover making them keep track of many file directories, sometimes with redundant information. An example of refinement would be to use web services to retrieve these files. We exemplify an even more helpful refinement – rather than reuse code, to reuse workflows that perform this task: we retrieved from the my-

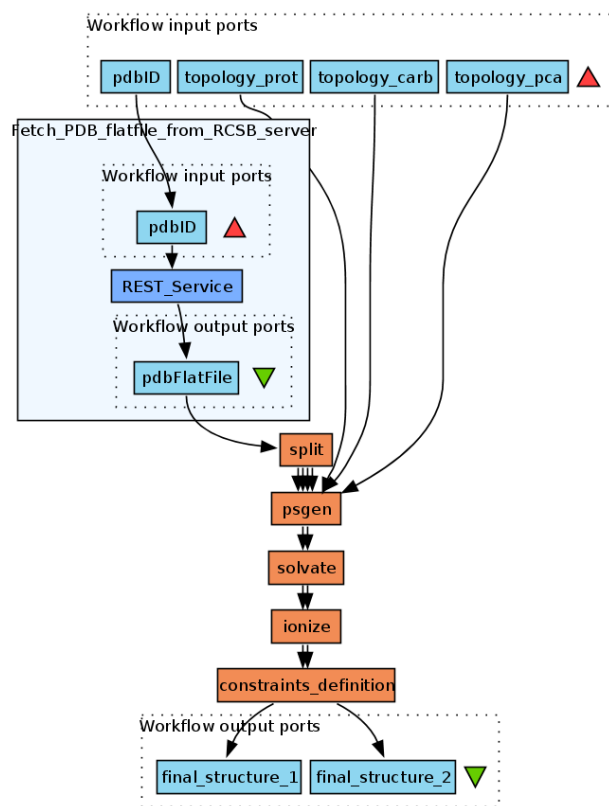


Fig. 7. Refined workflow of our MD case study, extracted from [14].

Experiment repository⁵ a small workflow that fetches a protein structure on Protein Data Bank (PDB) from the RCSB PDB archive⁶. The reused myExperiment workflow was inserted in the beginning of our original workflow (see figure 7 for the workflow variant We_1).

Here, the *initial_structure* input parameter of figure 6 (the local PDB file) was replaced by the sub-workflow within the light blue box, copied from the myExperiment workflow repositories. This new (sub)workflow downloads the protein file from the web using a web service (whereas the original code used a local protein file). Similarly, in the life sciences, scientists can invoke web services or reuse data sets listed on portals such as Biocatlogue⁷, which provides a curated catalogue of Web services, and Biotools⁸, which is a tools and data services registry.

During the conversion process, additional activities must be performed. First, it is critical to have a quality

⁵<http://www.myexperiment.org>

⁶<http://www.rcsb.org/pdb/>

⁷<https://www.biocatlogue.org/>

⁸<https://bio.tools>

1 check where the scientist explicitly assesses the work-
 2 flow activities and data flow, comparing them to what
 3 was executed by the script. Hence, throughout the pro-
 4 cess of workflow creation and modification, the scien-
 5 tist should provide annotations describing it (i.e. activi-
 6 ties and ports), and potentially the resources it utilizes.
 7 Part of these annotations can be migrated to the exe-
 8 cutable workflow taking advantage of the YesWork-
 9 flow tags - e.g., `@desc` used in the script to describe its
 10 program blocks and ports. Most SWfMSs, moreover,
 11 provide an annotation interface, which can be taken ad-
 12 vantage of. Our work [24] partially describes some of
 13 those annotation tasks.

14 Second, provenance information is used by W2Share
 15 for several purposes. Besides experiment reproducibil-
 16 ity, it is recorded to capture the steps performed
 17 in the transformation from script to workflow. This
 18 uses a provenance model, which allows identifying
 19 the correspondence between workflow activiti(es) and
 20 script code, and reusable components/web services
 21 and script excerpts. The lineage of variants of the
 22 workflow should be stored, as well. It is important to
 23 inform to future users that the workflow was curated,
 24 and how this curation process occurred. Provenance
 25 capture is presented in section 5 of this paper.

26 Finally, W2Share assists the scientist to create a
 27 Workflow Research Object (WRO) that bundles the
 28 original script as well as other auxiliary resources
 29 obtained in the other steps of the methodology. The
 30 Workflow Research Object model [9, 16] allows scien-
 31 tists to aggregate resources and explicitly specify the
 32 relationship between these resources and the workflow
 33 in a machine-readable format using a suite of ontolo-
 34 gies.

35 The resulting WRO bundles a number of resources
 36 that promote the understanding, reproducibility and ul-
 37 timately the reuse of the workflows obtained through
 38 refinement. By including these resources, it is possi-
 39 ble for scientists not only to understand how the ex-
 40 periment was conducted, but also its context. More-
 41 over, curators can also bundle additional documents
 42 that may help scientists understand the WRO, e.g.,
 43 technical reports and published papers.

44 We use the RO Manager tool⁹ to create the WRO
 45 bundle file. The bundle for this case study is available
 46 in [25].

47 However, it is not enough to create such research
 48 objects; they must be made available to the scientific
 49

1 community in a user-friendly manner, so that not only
 2 machines, but also scientists can select the most ap-
 3 propriate ones. A possible solution is to make them
 4 available by depositing them in a Research Object Por-
 5 tal such as W2Share¹⁰, myExperiment and RO Hub¹¹
 6 which have an interface to search and navigate be-
 7 tween resources aggregated in a RO.
 8

5. Revisiting the Implementation of Step 1: Mapping Scripts into PROV-Compliant Machine-Readable Abstract Workflows

14 The first step of our methodology is the conversion
 15 of scripts into abstract workflows. One of our innova-
 16 tions is the creation of a new kind of abstract workflow
 17 for scripts – one that is ontology-based and, moreover,
 18 machine-readable. We call this a "machine-readable
 19 abstract workflow" (as opposed to the abstract work-
 20 flows described in the literature, which are usually
 21 structures devoid of any semantics).

22 This section explains how W2Share enables the
 23 transformation of a script S into the machine-readable
 24 abstract workflow Wa . The latter, in turn, is used
 25 to create and describe the corresponding executable
 26 workflow We , and its subsequent variants We_1 , We_2 ,
 27 etc.

28 Section 4.2 shows how a scientist can easily trans-
 29 form a script into an abstract workflow with help of
 30 the YesWorkflow suite of tools [5, 26]. However, these
 31 abstract workflows are not machine-readable. Indeed,
 32 YesWorkflow has two outputs – an image of a work-
 33 flow, and Datalog code that encodes the correspond-
 34 ing structure. This limits its interoperability with ap-
 35 proaches that use semantic technologies. Moreover,
 36 YesWorkflow's workflow representation, if considered
 37 apart from the originating script, does not allow ob-
 38 taining provenance information on how it was derived
 39 from the script.

40 Our solution is to transform script S into machine-
 41 readable abstract workflows Wa in a three-stage pro-
 42 cess. First, we use YesWorkflow to extract the work-
 43 flow topology from S . Next, we transform this struc-
 44 ture into an ontology-based structure using a workflow
 45 specification ontology. Finally, we add provenance in-
 46 formation to link Wa back to S (thereby also support-
 47 ing traceability of the conversion process).
 48

51 ⁹<https://github.com/wf4ever/ro-manager>

¹⁰<https://w3id.org/w2share>

¹¹<http://www.rohub.org/>

Table 1
Mapping of YesWorkflow tags to workflow ontologies

Tag	Class	Property
@begin	wfdesc:Workflow	wfdesc:hasSubProcess
	wfdesc:Process	rdfs:label
	wf4ever:Script	
@desc	–	dct:description
@in	wfdesc:Input	wfdesc:hasInput; rdfs:label
@param		
@out	wfdesc:Output	wfdesc:hasOutput; rdfs:label
	wfdesc:DataLink	wfdesc:hasDataLink
		wfdesc:hasSource
		wfdesc:hasSink
@as	–	rdfs:label

In more detail, the first stage creates the YesWorkflow abstract representation. In the second stage, we use ontologies to transform this representation into a semantic one. This is achieved by mapping YesWorkflow tags to workflow entities that are semantically defined via wfdesc [9], a workflow specification ontology from the Research Object suite of ontologies [16], and other additional ontologies, e.g., Wf4ever¹², RDF Schema and Dublin Core¹³. Finally, in the third stage, we process tags to insert provenance information, i.e., we create an additional layer of provenance over the abstract semantic workflow.

Table 1 summarizes the second stage, showing how YesWorkflow tags are mapped to classes and properties of ontologies. Here, we use the following name spaces: ro for Research Object, wfdesc for the wfdesc Ontology, wf4ever for the Wf4ever Schema, dct for the Dublin Core terms, rdfs for RDF schema and prov for the PROV ontology.

Figure 8 shows an excerpt of the second stage. On the left side of this figure we have a script using YesWorkflow tags. The right side shows the RDF triples that correspond to these tags. Numbers ①, ②, ③, and ④ connect both sides of the figure. For instance, on the left side, ① has annotation @begin md_setup, which is mapped to the class wfdesc:Workflow and the property rdfs:label with value md_setup. While @desc setup of a MD simulation is mapped to property dc:description with value setup of a MD simulation.

¹²<https://w3id.org/ro/wf4ever>

¹³<http://dublincore.org/>

On the left side, ② originates the input and output RDF triples on the right side. The input and output ports of the workflow use classes wfdesc:Input and wfdesc:Output, and properties rdfs:label and dc:description. Nested @begin tag ③ (left side) are mapped to the wfdesc:hasSubProcess property and the wfdesc:Process class, specifying an activity of the workflow. The mapping also uses the wfdesc:DataLink class, as well as the properties wfdesc:hasSource and wfdesc:hasSink, identifying a link between two ports in the workflow.

At the third stage of the transformation of S to Wa , provenance information is added to the triples code. Provenance semantics are provided by using the PROV ontology [15]. Each abstract workflow element is defined as a prov:Entity. Again, in Figure 8, in ④, the script filename is mapped to triples defining the script resource <resources/script.sh> as a wf4ever:Script. This specific mapping is independent of the use of any YesWorkflow tag. The property prov:wasDerivedFrom is created with value <resources/script.sh>, identifying from which script that workflow was derived, since an experiment may have more than one script file. The script code committed within a block in ④ (left side), originates the identification of the text position in the script code using properties and classes such as prov:wasDerivedFrom from PROV and oa:TextPositionSelector, oa:start, and oa:end from the Web Annotation Ontology, to delimit this code. Another useful provenance information added at this stage is by whom and when the transformation was performed.

Summing up, this section described W2Share's process to generate a machine-readable (semantic) abstract representation of a workflow in which workflow blocks are linked back to the original script block, and script resources are duly semantically annotated. Annotations also indicate the agents responsible for the script-to-workflow conversion. This helps reproducibility by documenting the conversion process. This also helps reuse of workflow specifications.

Listing A.2 in A shows an excerpt of Wa and Listing A.3 in A shows the provenance information generated by the transformation. These listings use the RDF Turtle format¹⁴ and are results of the three transformation stages from S to Wa using our MD case study.

¹⁴<http://www.w3.org/TR/turtle/>

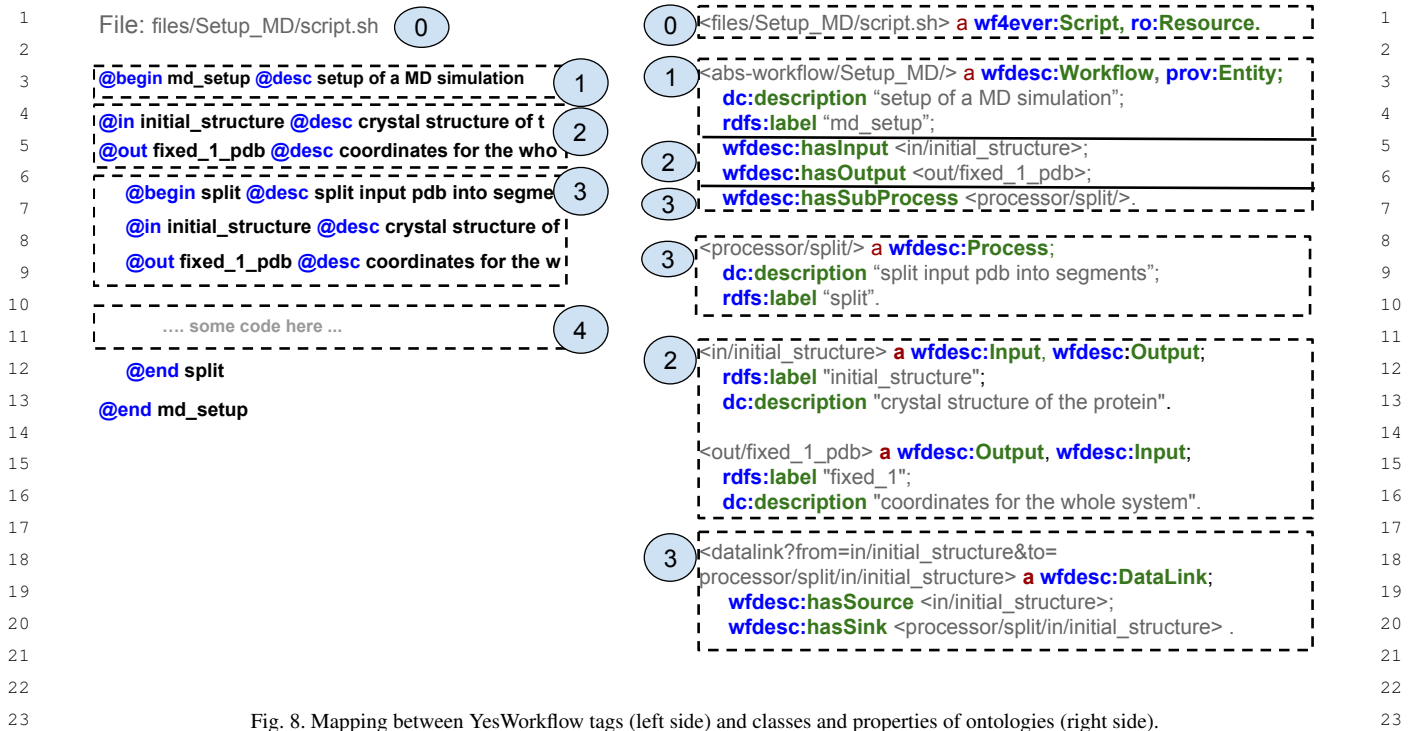


Fig. 8. Mapping between YesWorkflow tags (left side) and classes and properties of ontologies (right side).

6. Revisiting Step 2: (Semi-)Automatically Transforming Abstract Workflows into Executable Workflows

After creating Wa , a machine-readable abstract workflow, the next step is to create an executable workflow We , which corresponds to Step 2 of our methodology. We here show how this can be done automatically, in the best case and semi-automatically in the other cases, mapping Wa elements to elements that can be executed in a SWfMS, e.g., Taverna. This step was outlined in a previous work [24]; here we provide a detailed description.

As presented in Section 5, our machine-readable abstract workflow Wa describes the activities encoded by the script code. In our provenance layer, entities point back to the corresponding script code block. The Wa to We conversion process will now take advantage of this provenance information. At the end, We and S share the same abstract workflow Wa to describe their activities in a higher level.

During the creation of We , the original script S code may be manually changed (e.g., to allow appropriate workflow execution in the chosen SWfMS). So, the scientist must be aware of potential issues caused by these modifications. Some changes in the code can be performed automatically, e.g., library imports. Others

might need manual intervention, such as changing a reference to an absolute path to a file to obtain the file from a workflow port. By identifying these manual changes in the workflow implementation, experts can describe the reason behind the changes, which helps documenting the conversion process.

W2Share's machine-readable abstract workflows allow linking We to S , thus enabling questions related to the sequence of transformation steps that led to the production of an executable workflow. Examples include "which script block originated a specific activity in this workflow?", or "which workflow activities do not have exactly the same code as the script code that originated it?". The latter question would use the current implementation of the workflow activity and a simple comparison with the original script block code. In addition, scientists can use annotations regarding the reason behind the changes to foster the understanding of the process.

Listing A.3 in A shows an excerpt of PROV statements in RDF Turtle format to allow tracing back elements of the executable workflow We to the elements of script S through Wa .

7. Evaluation

7.1. Overview

To validate our proposal, we adopt the notion of *competency queries*, e.g., as defined by [27] "A competency query is a query that is useful for the community at hand, e.g. for a human member (e.g. a scientist), or for building applications for that domain. Therefore, a list of such queries can sketch the desired scope and the desired structuring of the information."

Our competency queries show W2Share's ability to answer questions about a workflow's lineage thanks to prospective provenance generated during the script-to-workflow process. Questions about workflow executions are answered thanks to retrospective provenance obtained from the SWfMS during workflow execution. All these resources are bundled in WROs.

The queries proposed in this section should help scientists to understand and explore the conversion process and consequently assess the quality and establish trust in this process. To achieve this goal, the queries return prospective and retrospective provenance information. Examples of the prospective view include, e.g, how the workflow was created from the script, and who created the workflow. Retrospective views include associating workflow results to the script which originated the workflow.

Given as input S , Wa , We , We_1 , ... We_n , we consider the following kinds of queries for prospective provenance:

1. tracking elements: activities, data, data flows in We back to S ;
2. metadata: information describing elements of script, workflows, and agents;
3. provenance of a given data source (before execution).

We consider the following kinds of queries for retrospective provenance:

1. establishing trust: comparison of workflow and script results, and comparison of workflow variant results;
2. tracking elements: link elements derived from S to traces.

Here, we consider that these queries are important to help the scientists to establish trust and assess the quality of the conversion by comparing the workflow results to the script results.

7.2. Executing Queries

Here, we specify competency questions associated to the requirements of Section 2. To each question, we provide a specific SPARQL query, which we evaluate against the contents of the WRO generated by our case study.

We point out that these queries do not inspect the scripts or executable workflows. Rather, queries are processed against machine-readable abstract workflows. All these prospective provenance representations use *wfdesc*, PROV and Web Annotation ontologies. The retrospective provenance representation uses *wfprov* ontology which is part of the Research Object suite of ontologies [16]. The competency questions, the SPARQL queries and results, and the RDF statements representing our case study can be found online in [28].

We are interested in the following competency queries to address the requirements:

1. Retrieving information about the abstract workflow Wa derived from S (i.e., processing units and their dependencies) – addressing requirement R1.
2. Retrieving information about workflow We derived from S – addressing requirement R2.
3. Retrieving information about workflow variants derived from We (i.e., We_1 , We_2 , ... We_n) – addressing requirement R3.
4. Retrieving lineage information associating We and script elements (i.e., input, outputs, activities and data links) – addressing requirement R4.
5. Retrieving metadata about the conversion process (e.g, the person who created a given executable workflow) – addressing requirement R4.
6. Retrieving information tracking workflow execution traces to script blocks – addressing requirement R4.
7. Retrieving the resources available in the WRO associated to an experiment – addressing requirement R5.

We now present a description of each query, why it is relevant, the kind of situations for which it would be needed, their translation into SPARQL¹⁵, and the results obtained by evaluating them.

These queries are run against the RDF statements in Listings A.2, A.3, A.4, and A.5 in Appendix A.

¹⁵<https://www.w3.org/TR/rdf-sparql-query/>

Query 1 Retrieving information about abstract workflow Wa derived from S . This query is responsible for identifying, given S , the processing units and their dependencies which compose Wa . It is useful, for example, for scientists to understand the data analysis carried out by the experiment.

Listing 1 shows the SPARQL query that can be used for answering this query. We use the WRO URI as base for the queries: <https://w3id.org/w2share/wro/md-setup/>.

Listing 1: Query 1 translated into a SPARQL

```

1 select ?abs ?winput ?woutput ?process ?pinput
2   ↳ ?poutput
3 where {
4   ?abs prov:wasDerivedFrom
5     ↳ <files/Setup_MD/script.sh> ;
6     wfdesc:hasInput ?winput ;
7     wfdesc:hasOutput ?woutput ;
8     wfdesc:hasSubProcess ?process .
9   ?process wfdesc:hasInput ?pinput ;
10    wfdesc:hasOutput ?poutput .
11 }

```

Table 2 shows the results obtained by evaluating the query. The results point out that Wa `<abs-workflow/Setup_MD/>` is the abstract workflow derived from `<files/Setup_MD/script.sh>`. Wa has input `<abs-workflow/Setup_MD/in/structure_pdb>` and output `<abs-workflow/Setup_MD/out/fixed_1_pdb>`. Also, it has `<abs-workflow/Setup_MD/processor/split>` as one of its activities, which has input `<abs-workflow/Setup_MD/processor/split/in/structure_pdb>` and output `<abs-workflow/Setup_MD/processor/split/out/fixed_1_pdb>`.

Table 2
Result of evaluating query 1

Variable	Value
abs	<abs-workflow/Setup_MD>
winput	<abs-workflow/Setup_MD/in/structure_pdb>
woutput	<abs-workflow/Setup_MD/out/fixed_1_pdb>
processor	<abs-workflow/Setup_MD/processor/split>
pinput	<abs-workflow/Setup_MD/processor/split/in/initial_structure>
poutput	<abs-workflow/Setup_MD/processor/split/out/cbh1_pdb>

Query 2: Retrieving information about the executable Workflow derived from S . This query is re-

sponsible for identifying which executable workflow is derived from S . This is useful for scientists interested in executing or reusing pieces of this workflow.

Listing 2: Query 2 translated into a SPARQL

```

1 select ?workflow
2 where {
3   ?workflow a wfdesc:Workflow ;
4   prov:wasDerivedFrom
5     ↳ <files/Setup_MD/script.sh> .
6 }

```

Listing 2 shows the SPARQL code for this query. Table 3 shows the results obtained by evaluating the query. The results point out that `<workflow/Setup_MD>` is derived from `<files/Setup_MD/script.sh>`.

Table 3
Result of evaluating query 2

Variable	Value
Workflow	<workflow/Setup_MD>

To have a deep understanding of the differences between the implementation of the workflow We and the script S , it would be necessary to compare the specification of both implementations, or to have annotations describing the rationale of these changes. However, this is outside the scope of this paper.

Query 3: Retrieving information about workflow variants derived from We . This query is responsible for identifying, given an executable workflow We , which workflows are derived from it. It is useful, for example, for scientists to find workflow variants to run, reuse or compare their results. This query is also useful when a scientist updates an executable workflow, and needs to propagate this update to the workflow variants.

Listing 3 shows the SPARQL code that can be used for answering this query.

Listing 3: Query 3 translated into a SPARQL

```

1 select ?variant
2 where {
3   ?variant prov:wasDerivedFrom
4     ↳ <workflow/Setup_MD/> .
5 }

```

Table 4 shows the results obtained by evaluating the query. The results point out `<workflow/Setup_MD>` originated the variant `<workflow/Setup_MD/variant>`.

Table 4
Result of evaluating query 3

Variable	Value
variant	<code><workflow/Setup_MD/variant></code>

Again, to have a deep understanding of the differences between the implementations of the workflow We and its variants We_1, We_2, \dots, We_n , it would be necessary to compare the specification of both implementations, or to have annotations describing the rationale of these changes. Again, this is outside the scope of this paper.

Query 4: Retrieving lineage information associating We and script elements.

This query is responsible for identifying, given an executable workflow We , which script blocks originated each workflow activity. It is useful, for example, for scientists to compare the executable workflow and script implementations.

We use the Web Annotation Ontology element `oa:TextPositionSelector` and its properties `oa:start` and `oa:end` to delimit the textual position of blocks in the script code. Listing 4 shows the SPARQL code that can be used for answering this query.

Listing 4: Query 4 translated into a SPARQL

```

1 select ?abs ?process ?start ?end
2 where {
3   ?abs prov:wasDerivedFrom
4     ↪ <files/Setup_MD/script.sh> ;
5   wfdesc:hasSubProcess ?process .
6   ?process prov:wasDerivedFrom ?code .
7   ?code oa:start ?start ;
8     oa:end ?end .
}

```

Table 5 shows the results obtained by evaluating the query. The results point out that `<workflow/Setup_MD/split>` was derived from `<abs-workflow/Setup_MD/split>`, which is defined in the text position from 1644 to 1786 in S .

Query 5: Retrieving metadata about the conversion process. This query is responsible for retrieving metadata regarding each step of the script-to-workflow conversion process (e.g., who was the person in charge

Table 5
Result of evaluating query 4

Variable	Value
workflow_process	<code><workflow/Setup_MD/processor/split></code>
script_process	<code><abs-workflow/Setup_MD/processor/split></code>
block_start	1644
block_end	1786

of annotating a script S to create the abstract workflow).

This query is useful helping to establish trust in the conversion process. Listing 5 shows the corresponding SPARQL code. Table 6 shows the result of the query, which points out Lucas Carvalho is the curator associated with the creation of `<abs-workflow/Setup_MD/>` for S .

Listing 5: Query 5 translated into a SPARQL

```

1 select distinct ?curator
2 where {
3   <abs-workflow/Setup_MD/>
4     ↪ prov:wasAttributedTo ?agent .
5   ?agent foaf:name ?curator .
}

```

Table 6
Result of evaluating query 5

Variable	Value
curator	Lucas Carvalho

Query 6: Retrieving information tracking workflow execution traces of We to script blocks. This query is responsible for identifying, given a workflow execution trace, the original script blocks associated with it. This query is useful when a scientist wants to retrieve workflow executions and compare them with script executions.

Listing 6 shows the SPARQL code for this query. Table 7 shows the result of the query, which points out that `<workflow/processor/split/>` was derived from the script block `<abs-workflow/processor/split>`, used as input `<>` and produced as output `<data/4e0alf-fc0f/output/bg1c.pdb>`.

Listing 6: Query 6 translated into a SPARQL

```

1 SELECT DISTINCT ?workflow ?workflowRun
2   ↪ ?process ?output ?input
3 WHERE {

```

```

1 3   ?workflow prov:wasDerivedFrom
2     ↪ <files/Setup_MD/script.sh> .
3 4   ?workflow wfdesc:hasSubProcess ?process .
4 5   ?processRun wfprov:describedByProcess
5     ↪ ?process ;
6     prov:used ?input ;
7     wfprov:wasPartOfWorkflowRun
8     ↪ ?workflowRun .
8 8   ?output prov:wasGeneratedBy ?processRun .
9 9 }

```

Table 7
Result of evaluating query 6

Variable	Value
workflow	<workflow/Setup_MD/>
workflowRun	<run/4e0a1f-fc0f/>
process	<workflow/Setup_MD/processor/split/>
input	<data/4e0a1f-fc0f/input/structure.pdb>
output	<data/4e0a1f-fc0f/output/blgc.pdb>

Query 7: Retrieving the resources available in a WRO. This query is responsible for identifying, given a WRO, which resources are aggregated by it. This query is useful to identify which resources to reuse, for example.

Table 7 shows the result of the query, which points out that the WRO aggregates the resources `script.sh`, `executable-workflow.t2flow`, `refined-workflow.t2flow`, and their inputs and output files (resources) aggregated in the WRO.

Listing 7: Query 7 translated into a SPARQL

```

1 select distinct ?resource ?type
2 where {
3   <> ore:aggregates ?resource .
4   ?resource a ?type .
5   FILTER(ro:Resource != ?type)
6 }

```

Table 8
Result of evaluating query 7

resource	type
<workflow/executable-workflow.t2flow>	wf4ever:Workflow
<workflow/refined-workflow.t2flow>	wf4ever:Workflow
<files/script.sh>	wf4ever:Script
<workflowrun.prov.ttl>	wfdesc:WorkflowRun
<data/4e0a1f-fc0f/input/structure.pdb>	wf4ever:Dataset
<data/4e0a1f-fc0f/output/blgc.pdb>	wf4ever:Dataset

8. Related Work

Here, we provide a comparison of our work with research on script-to-workflow conversion and traceability. Additional related work was already commented on throughout the text.

Scripts are usually difficult to understand, reuse, or reproduce by people other than the original implementers. In previous work [14], we described a preliminary manual transformation from script to executable and verifiable workflow. We adopted YesWorkflow to generate the abstract workflow visualizations before creating the executable workflow. YesWorkflow was developed by McPhilips *et al.* [5, 26]; it is a script language-independent environment for extracting a workflow-like view that depicts the main components that compose a script and their data dependencies based on comments that annotate the script.

At the time of writing this paper, the YesWorkflow group is working towards supporting the exportation, in RDF, of an abstract workflow representation that is PROV-compliant, using the ProvONE model [29]. We, instead, use *wfdesc* and *wfprov* ontologies as PROV extensions because we are targeting Research Objects, and also the Taverna SWfMS. Moreover, our provenance information is used to link the abstract workflow back to the script.

There are a few other approaches to construct executable workflows from scripts. For instance, [10] uses the abstract syntax tree (AST) created from source code to map the script elements into workflow structures. Our approach differs from this in that we reuse parts of the script code to create the workflow activities.

The work of [11] migrates script-based experiments from a local High Performance Computer (HPC) cluster to workflows on a cloud computing infrastructure. Their requirements include traceability of the workflow results to meet reproducibility, one of their reasons to migrate to a SWfMS. One of the differences to our approach is that our methodology is more generic, in the sense that we do not focus on HPC scripts nor on cloud computing environments, which require specific kinds of scripts. Also, we do not consider any specific approach to meet the challenge of converting scripts with control-flow constructs into data-flow patterns, which is addressed in both [11] and [10].

In [12], the authors present an approach to convert electronic notebooks into workflows. Their approach goes directly from notebook code to executable workflow. It is based on a set of guidelines that recom-

1 mend changes to the notebook structure to facilitate
2 the capture of the dataflow encoded in the notebook
3 and enable its conversion into an executable workflow.
4 We, instead, go from script-to-abstract and abstract-
5 to-executable steps, thereby clearly separating abstract
6 specification from code. This helps documentation, un-
7 derstandability and reuse. At the end, their conversion
8 is performed by NiW, a tool that semi-automatically
9 creates the workflow structure based on the notebook's
10 code. Similar to our approach, after running their tool,
11 a scientist may need to improve the corresponding
12 workflow implementation due to differences in the en-
13 vironments. Our annotation and abstract workflows
14 guide the scientists to identify the main processing
15 units and dataflow in the script.

16 In [24], we presented the first implementation of a
17 web prototype for W2Share that integrates the tools
18 used to convert scripts into WRO. In that paper, we
19 also address quality checking of the conversion pro-
20 cess. However, we do not address tracking issues dur-
21 ing the conversion process. This is a recent result that
22 is being reported here.

23 The work of [30] presents an approach to track
24 changes in workflows to capture the evolution of work-
25 flows and allow the comparison of results and struc-
26 ture between different versions. We, instead, focus
27 on issues related to the traceability of the script-to-
28 workflow conversion process, which enables relating
29 workflow elements back to the original script blocks,
30 comparing differences between script and workflow
31 implementation, and comparing differences between
32 script and workflow results.

33 Another difference between ours and other script-to-
34 workflow approaches is that ours use of WROs. These
35 objects bundle the executable workflow, but also the
36 workflow specification and auxiliary resources, as well
37 as workflow runs and data used in these runs. Thus,
38 one single object (the WRO) is needed to support ex-
39 periment reproducibility, reuse, and checking of exper-
40 iments in a transparent way.

41 9. Conclusions and Ongoing Work

42
43
44
45 This paper is a step towards fully reproducible re-
46 search. It presented W2Share, a computation frame-
47 work that supports a (script-to-reproducible research)
48 methodology. The methodology, implemented in W2-
49 Share via a suite of tools, guides scientists in a prin-
50 cipled manner to transform scripts into reproducible
51 and reusable research objects. W2Share addresses an

1 important issue in the area of provenance of scientific
2 experiments modeled as scripts – that of providing an
3 executable and understandable provenance representa-
4 tion of domain script runs. We point out that prove-
5 nance is not just metadata for others: "provenance-
6 for-self" queries can be used by researchers to better
7 understand experiments, and to speed up the conver-
8 sion process. Thus, there is a need for support to hy-
9 brid provenance queries (i.e., involving both
10 prospective and retrospective queries).

11 Our ontology-based approach to generate machine-
12 readable abstract workflows is also useful for querying
13 purposes (e.g., traceability). It also allows associating
14 the executable workflow, represented using the wfdesc
15 ontology, and provenance information, represented us-
16 ing the wfprov ontology, in ontology-based queries.

17 W2Share was elaborated based on requirements
18 that we elicited given our experience and collabora-
19 tions with scientists who use scripts in their simula-
20 tions. Moreover, it enables traceability of the script-
21 to-workflow process, thereby establishing trust in this
22 process. The approach was showcased via a real
23 world use case from Molecular Dynamics. We showed
24 through competency questions that W2Share success-
25 fully meets those requirements. The competency ques-
26 tions and the case studies are additional contribu-
27 tions of our work. An initial implementation of our
28 methodology is described in [24] and available at
29 <https://w3id.org/w2share>.

30 Our ongoing and future work include promoting the
31 use of the conversion process in an e-Science infras-
32 tructure, investigating further real use cases with the
33 objective of extending it to fit (new) user requirements
34 and other script environments. Also, we plan to evalu-
35 ate the cost effectiveness of our proposal, in particu-
36 lar since in some cases it may require extensive in-
37 volvement of scientists. Last but not least, we do not
38 consider versioning. Thus, yet another future direction
39 would be to provide support to such version control
40 when refining executable workflows, for instance, by
41 considering an Ontology of Research Object Evolu-
42 tion.

43 Acknowledgments

44
45 Work partially financed by Sao Paulo Science Founda-
46 tion (FAPESP) under grants #2014/23861-4, #2017/
47 03570-3, FAPESP/CEPID CCES under grant #2013/
48 08293-7, and individual grants from CNPq. We thank
49 Prof. Munir Skaf and his group from the Institute of
50
51

Chemistry at Unicamp for making their scripts and data available and for their valuable feedback in the molecular dynamics case study.

References

- [1] E. Deelman, T. Peterka, I. Altintas, C.D. Carothers, K.K. van Dam, K. Moreland, M. Parashar, L. Ramakrishnan, M. Taufer and J. Vetter, The future of scientific workflows, *The International Journal of High Performance Computing Applications (IJHPCA)* **32**(1) (2018), 159–175.
- [2] S. Cohen-Boulakia, K. Belhajjame, O. Collin, J. Chopard, C. Froidevaux, A. Gaignard, K. Hinsén, P. Larmande, Y. Le Bras, F. Lemoine et al., Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities, *Future Generation Computer Systems* **75** (2017), 284–298.
- [3] T. Kluyver, B. Ragan-Kelley, F. Pérez, B.E. Granger, M. Bussonnier, J. Frederic, K. Kelley, J.B. Hamrick, J. Grout, S. Corlay et al., Jupyter Notebooks—a publishing format for reproducible computational workflows., in: *ELPUB*, 2016, pp. 87–90.
- [4] F. Chirigati, R. Rampin, D.E. Shasha and J. Freire, ReproZip: Computational Reproducibility With Ease., in: *SIGMOD Conference*, ACM, 2016, pp. 2085–2088. ISBN 978-1-4503-3531-7.
- [5] T. McPhillips, T. Song, T. Kolisnik, S. Aulenbach, K. Belhajjame, R.K. Bocinsky, Y. Cao, J. Cheney, F. Chirigati, S. Dey et al., YesWorkflow: A User-Oriented, Language-Independent Tool for Recovering Workflow Information from Scripts, *International Journal of Digital Curation* **10**(1) (2015), 298–313.
- [6] L. Murta, V. Braganholo, F. Chirigati, D. Koop and J. Freire, noworkflow: Capturing and analyzing provenance of scripts, in: *Provenance and Annotation of Data and Processes*, Springer, 2014, pp. 71–83.
- [7] S. Dey, K. Belhajjame, D. Koop, M. Raul and B. Ludäscher, Linking prospective and retrospective provenance in scripts, in: *7th USENIX Workshop on the Theory and Practice of Provenance (TaPP 15)*, 2015.
- [8] J. Liu, E. Pacitti, P. Valduriez and M. Mattoso, A Survey of Data-Intensive Scientific Workflow Management, *Journal of Grid Computing* **13**(4) (2015), 457–493, ISSN 1572-9184.
- [9] K. Belhajjame, O. Corcho, D. Garijo, J. Zhao, P. Missier, D. Newman, R. Palma, S. Bechhofer, E. García Cuesta, J.M. Gómez-Pérez et al., Workflow-Centric Research Objects: First Class Citizens in Scholarly Discourse, in: *Proceedings of Workshop on the Semantic Publishing, (SePublica 2012)*, 2012.
- [10] M. Baranowski, A. Belloum, M. Bubak and M. Malawski, Constructing workflows from script applications, *Scientific Programming* **20**(4) (2012), 359–377.
- [11] J. Cala, Y. Xu, E.A. Wijaya and P. Missier, From scripted HPC-based NGS pipelines to workflows on the cloud, in: *Cluster, Cloud and Grid Computing (CCGrid), 2014 14th IEEE/ACM International Symposium on*, IEEE, 2014, pp. 694–700.
- [12] L.A.M.C. Carvalho, R. Wang, D. Garijo and Y. Gil, NiW: Converting Notebooks into Workflows to Capture Dataflow and Provenance, in: *2017 Workshop on Capturing Scientific Knowledge (SciKnow), held in conjunction with the ACM International Conference on Knowledge Capture (K-CAP), December 4-6, Austin, TX, USA, 2017*, pp. 1–8.
- [13] J. Freire, D. Koop, E. Santos and C.T. Silva, Provenance for computational tasks: A survey, *Computing in Science & Engineering* **10**(3) (2008), 11–21.
- [14] L.A.M.C. Carvalho, K. Belhajjame and C.B. Medeiros, Converting Scripts into Reproducible Workflow Research Objects, in: *Proc. of the IEEE 12th Int. Conf. on eScience, October 23-26, IEEE, Baltimore, MD, USA, 2016*, pp. 71–80.
- [15] T. Lebo, S. Sahoo, D. McGuinness, K. Belhajjame, J. Cheney, D. Corsar, D. Garijo, S. Soiland-Reyes, S. Zednik and J. Zhao, PROV-O: the PROV ontology. W3C Recommendation, *World Wide Web Consortium* (2013). <https://www.w3.org/TR/prov-o/>.
- [16] K. Belhajjame, J. Zhao, D. Garijo, M. Gamble, K. Hettné, R. Palma, E. Mina, O. Corcho, J.M. Gómez-Pérez, S. Bechhofer et al., Using a suite of ontologies for preserving workflow-centric research objects, *Web Semantics: Science, Services and Agents on the World Wide Web* **32** (2015), 16–42.
- [17] P. Missier, S. Woodman, H. Hiden and P. Watson, Provenance and data differencing for workflow reproducibility analysis, *Concurrency and Computation: Practice and Experience* **28**(4) (2016), 995–1015, ISSN 1532-0634.
- [18] R. Souza, V. Silva, A.L. Coutinho, P. Valduriez and M. Mattoso, Data reduction in scientific workflows using provenance monitoring and user steering, *Future Generation Computer Systems* (2017). doi:10.1016/j.future.2017.11.028.
- [19] M. Mattoso, J. Dias, K.A. Ocaña, E. Ogasawara, F. Costa, F. Horta, V. Silva and D. de Oliveira, Dynamic steering of HPC scientific workflows: A survey, *Future Generation Computer Systems* **46** (2015), 100–113.
- [20] L.A.M.C. Carvalho, B.T. Essawy, D. Garijo, C.B. Medeiros and Y. Gil, Requirements for Supporting the Iterative Exploration of Scientific Workflow Variants, in: *2017 Workshop on Capturing Scientific Knowledge (SciKnow), held in conjunction with the ACM International Conference on Knowledge Capture (K-CAP), 2017*, pp. 1–8.
- [21] L.A.M.C. Carvalho, D. Garijo, C.B. Medeiros and Y. Gil, Semantic Software Metadata for Workflow Exploration and Evolution, in: *Proc. of the IEEE 14th Int. Conf. on eScience, Oct 29-Nov 1, IEEE, Amsterdam, Netherlands, 2018*.
- [22] R.L. Silveira and M.S. Skaf, Molecular dynamics simulations of family 7 cellobiohydrolase mutants aimed at reducing product inhibition, *The Journal of Physical Chemistry B* **119**(29) (2014), 9295–9303.
- [23] K. Wolstencroft, R. Haines, D. Fellows, A. Williams, D. Withers, S. Owen, S. Soiland-Reyes, I. Dunlop, A. Nenadic, P. Fisher, J. Bhagat, K. Belhajjame, F. Bacall, A. Hardisty, A. Nieva de la Hidalga, M.P. Balcazar Vargas, S. Sufi and C. Goble, The Taverna workflow suite: designing and executing workflows of Web Services on the desktop, web or in the cloud, *Nucleic Acids Research* **41**(W1) (2013), 557–561.
- [24] L.A.M.C. Carvalho, J.E.G. Malaverri and C.B. Medeiros, Implementing W2Share: Supporting Reproducibility and Quality Assessment in eScience, in: *Proc. of the 11th Brazilian e-Science Workshop (BreSci), July 5-6, 2017, Brazilian Computer Society, Sao Paulo, Brazil, 2017*.
- [25] L.A.M.C. Carvalho and C.B. Medeiros, W2Share Case Study: Workflow Research Object (WRO), 2018. doi:10.5281/zenodo.1465897.

- [26] T. McPhillips, S. Bowers, K. Belhajjame and B. Ludäscher, Retrospective provenance without a runtime provenance recorder, in: *7th USENIX Workshop on the Theory and Practice of Provenance (TaPP 15)*, 2015.
- [27] Y. Tzitzikas, N. Minadakis, Y. Marketakis, P. Fafalios, C. Allocca, M. Mountantonakis and I. Zidianaki, MatWare: Constructing and Exploiting Domain Specific Warehouses by Aggregating Semantic Data, in: *The Semantic Web: Trends and Challenges: 11th International Conference, ESWC 2014, Anisaras, Crete, Greece, May 25-29, 2014. Proceedings*, V. Prestiti, C. d'Amato, F. Gandon, M. d'Aquino, S. Staab and A. Tor-dai, eds, Springer International Publishing, 2014, pp. 721–736.
- [28] L.A.M.C. Carvalho and C.B. Medeiros, W2Share Evaluation: Competency Questions, 2018. doi:10.5281/zenodo.1465893.
- [29] V. Cuevas-Vicentín, B. Ludäscher, P. Missier, K. Belhajjame, F. Chirigati, Y. Wei and B. Leinfelder, Provone: A prov extension data model for scientific workflow provenance, 2015. <http://purl.org/provone>.
- [30] J. Freire, C.T. Silva, S.P. Callahan, E. Santos, C.E. Scheidegger and H.T. Vo, Managing Rapidly-Evolving Scientific Workflows, in: *Provenance and Annotation of Data*, L. Moreau and I. Foster, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2006, pp. 10–18. ISBN 978-3-540-46303-0.

Appendix A. Molecular Dynamics Case Study

Resources of the case study from Molecular Dynamics.

Listing A.1: Excerpt of an annotated MD script using YesWorkflow tags.

```

1 #!/bin/bash
2
3 # @BEGIN setup @DESC setup of a MD simulation
4 # @PARAM directory_path @AS directory
5 # @IN initial_structure @DESC PDB: 8CEL
6 # @URI file:{directory}/structure.pdb
7 # @IN topology_prot
8 # @URI file:top_all122_prot.rtf
9 # @IN topology_carb
10 # @URI file:top_all122_prot.rtf
11 # @OUT gh5_psf @AS final_structure_psf
12 # @URI file:{directory}/gh5.psf
13 # @OUT gh5_pdb @AS final_structure_pdb
14 # @URI file:{directory}/gh5.pdb
15
16 # @BEGIN split
17 # @IN initial_structure
18 # @URI file:structure.pdb
19 # @IN directory_path @AS directory
20 # @OUT protein_pdb
21 # @URI file:{directory}/protein.pdb
22 # @OUT bglc_pdb
23 # @URI file:{directory}/bglc.pdb
24 # @OUT water_pdb
25 # @URI file:{directory}/water.pdb

```

```

26 structure = $directory_path"/structure.pdb"
27 protein = $directory_path"/protein.pdb"
28 water = $directory_path"/water.pdb"
29 bglc = $directory_path"/bglc.pdb"
30 egrep -v '(TIP3|BGLC)' $structure > $protein
31 grep TIP3 $structure > $water
32 grep BGLC $structure > $bglc
33 # @END split
34
35 # @BEGIN psfgen @DESC generate the PSF file
36 # @PARAM topology_prot
37 # @URI file:top_all122_prot.rtf
38 # @PARAM topology_carb
39 # @URI file:top_all136_carb.rtf
40 # @IN protein_pdb
41 # @URI file:protein.pdb
42 # @IN bglc_pdb
43 # @URI file:bglc.pdb
44 # @IN water_pdb
45 # @URI file:water.pdb
46 # @OUT hyd_pdb
47 # @URI file:hyd.pdb
48 # @OUT hyd_psf
49 # @URI file:hyd.psf
50
51 ... commands ...
52
53 # @END psfgen
54
55 # @BEGIN solvate
56 # @IN hyd_pdb
57 # @URI file:hyd.pdb
58 # @IN hyd_psf
59 # @URI file:hyd.psf
60 # @OUT wbox_pdb
61 # @URI file:wbox.pdb
62 # @OUT wbox_psf
63 # @URI file:wbox.psf
64 echo "
65 package require solvate
66 solvate hyd.psf hyd.pdb -rotate -t 16 -o wbox
67 exit
68 " > solv.tcl
69
70 vmd -dispdev text -e solv.tcl
71 rm solv.tcl
72 # @END solvate
73
74 # @BEGIN ionize
75 # @IN wbox_pdb
76 # @URI file:wbox.pdb
77 # @IN wbox_psf
78 # @URI file:wbox.psf
79 # @OUT gh5_pdb @AS final_structure_pdb
80 # @URI file:gh5.pdb
81 # @OUT gh5_psf @AS final_structure_psf
82 # @URI file:gh5.psf
83
84 ... commands ...
85
86 # @END ionize

```

```

1 87
2 88 # @END setup
3
4
5 Listing A.2: Excerpt of specification of Wa, the result
6 of transforming script S into an equivalent machine-
7 readable abstract workflow.
8
9
10 @base <https://w3id.org/w2share/wro/md-setup/
11   ↪ abs-workflow/Setup_MD/>.
12 @prefix dcterms: <http://purl.org/dc/terms/>.
13 @prefix wf4ever:
14   ↪ <http://purl.org/wf4ever/wf4ever#>.
15 @prefix oa:
16   ↪ <http://www.w3.org/ns/oa#>.
17 @prefix wfdesc:
18   ↪ <http://purl.org/w4ever/wfdesc#>.
19 @prefix prov:
20   ↪ <http://www.w3.org/ns/prov-o#>.
21 @prefix xsd:
22   ↪ <http://www.w3.org/2001/XMLSchema#>.
23 @prefix rdfs:
24   ↪ <http://www.w3.org/2000/01/rdf-schema#>.
25
26 <>
27
28 a wfdesc:Workflow, prov:Entity;
29 rdfs:label "setup"^^xsd:string;
30 wfdesc:hasSubProcess
31   ↪ <processor/split/>;
32 wfdesc:hasInput <in/structure_pdb>;
33 wfdesc:hasOutput <out/fixed_1_pdb> .
34
35 <in/initial_structure>
36 a wfdesc:Input, wfdesc:Output;
37 rdfs:label "structure_pdb"^^xsd:string;
38 dcterms:title "crystal structure of the
39   ↪ protein"^^xsd:string .
40
41 <out/fixed_1_pdb>
42 a wfdesc:Output, wfdesc:Input;
43 rdfs:label "fixed_1"^^xsd:string;
44 dcterms:title "coordinates for the whole
45   ↪ system (cbhl.pdb), indicating which
46   ↪ atoms should be kept fixed along the
47   ↪ simulation"^^xsd:string .
48
49 <datalink?from=in/initial_structure&to=
50   ↪ processor/split/in/initial_structure>
51 a wfdesc:DataLink;
52 wfdesc:hasSource <in/initial_structure>;
53 wfdesc:hasSink
54   ↪ <processor/split/in/initial_structure> .
55
56 <processor/split/>
57 a wfdesc:Process;
58 rdfs:label "split"^^xsd:string;
59 wfdesc:hasInput
60   ↪ <processor/split/in/initial_structure>;
61 wfdesc:hasOutput
62   ↪ <processor/split/out/cbhl_pdb> .

```

```

37
38 <processor/split/in/initial_structure>
39 a wfdesc:Input;
40 rdfs:label
41   ↪ "structure_pdb"^^xsd:string;
42 dcterms:description "crystal structure of
43   ↪ the protein"^^xsd:string .
44
45 <processor/split/out/cbhl_pdb>
46 a wfdesc:Output;
47 rdfs:label "cbhl_pdb"^^xsd:string;
48 dcterms:description "coordinates of the
49   ↪ protein atoms"^^xsd:string .

```

Listing A.3: Excerpt of PROV-statements describing the derivation of *S* to *Wa* to *We* to *We*₁.

```

1 @base
2   ↪ <https://w3id.org/w2share/wro/md-setup/>.
3 @prefix dcterms: <http://purl.org/dc/terms/>.
4 @prefix wf4ever:
5   ↪ <http://purl.org/wf4ever/wf4ever#>.
6 @prefix oa:
7   ↪ <http://www.w3.org/ns/oa#>.
8 @prefix wfdesc:
9   ↪ <http://purl.org/w4ever/wfdesc#>.
10 @prefix prov:
11   ↪ <http://www.w3.org/ns/prov-o#>.
12 @prefix xsd:
13   ↪ <http://www.w3.org/2001/XMLSchema#>.
14 @prefix rdfs:
15   ↪ <http://www.w3.org/2000/01/rdf-schema#>.
16 @prefix foaf: <http://xmlns.com/foaf/0.1/>.
17
18 <files/Setup_MD/script.sh> a wf4ever:Script.
19
20 <abs-workflow/Setup_MD/>
21   prov:wasDerivedFrom
22     ↪ <files/Setup_MD/script.sh> ;
23   prov:wasAttributedTo [
24     a prov:Agent ;
25     foaf:name "Lucas Carvalho" ] .
26
27 <abs-workflow/Setup_MD/processor/split/>
28   prov:wasDerivedFrom [
29     a prov:Entity, oa:TextPositionSelector;
30     oa:start "1644"^^xsd:Integer;
31     oa:end "1786"^^xsd:Integer;
32   ] .
33
34 <workflow/Setup_MD/> a wfdesc:Workflow,
35   ↪ prov:Entity ;
36   prov:wasDerivedFrom
37     ↪ <files/Setup_MD/script.sh> ;
38   prov:wasDerivedFrom
39     ↪ <abs-workflow/Setup_MD/> ;
40   wfdesc:hasSubProcess
41     ↪ <workflow/Setup_MD/processor/split/> .

```

```

1 32 <workflow/Setup_MD/processor/split/>
2 33   prov:wasDerivedFrom
3     ↪ <abs-workflow/Setup_MD/processor/split/>
4     ↪ .
5 34
6 35 <workflow/Setup_MD/variant> a
7     ↪ wfdesc:Workflow, prov:Entity ;
8 36   prov:wasDerivedFrom <workflow/Setup_MD/>
9     ↪ .

```

Listing A.4: Excerpt of workflow execution traces of We.

```

14
15 1 @base
16     ↪ <https://w3id.org/w2share/wro/md-setup/>
17     ↪ .
18 2 @prefix prov: <http://www.w3.org/ns/prov#> .
19 3 @prefix wfprov:
20     ↪ <http://purl.org/wf4ever/wfprov#> .
21 4 @prefix rdfs:
22     ↪ <http://www.w3.org/2000/01/rdf-schema#>
23     ↪ .
24 5 @prefix wfdesc:
25     ↪ <http://purl.org/wf4ever/wfdesc#> .
26 6 @prefix tavernaprov:
27     ↪ <http://ns.taverna.org.uk/2012/
28     ↪ tavernaprov/> .
29 7 @prefix owl:
30     ↪ <http://www.w3.org/2002/07/owl#> .
31 8 @prefix xsd:
32     ↪ <http://www.w3.org/2001/XMLSchema#> .
33 9 @prefix rdf:
34     ↪ <http://www.w3.org/1999/02/22-rdf-
35     ↪ syntax-ns#> .
36 10 @prefix dct: <http://purl.org/dc/terms/> .
37 11
38 12 <run/e0fa2f25-0755/>
39 13   rdf:type wfprov:WorkflowRun ;
40 14   dct:hasPart
41     ↪ <run/e0fa2f25-0755/process/f0a0bd65-78d3/>
42     ↪ ;
43 15   wfprov:describedByWorkflow
44     ↪ <workflow/Setup_MD/> ;
45 16   prov:used
46     ↪ <data/5c65c151-0333/ref/61f8795e-e650> ;
47 17   dct:hasPart
48     ↪ <run/e0fa2f25-0755/process/c06ff05e-eceb/>
49     ↪ ;
50 18   prov:endedAtTime
51     ↪ "2016-06-16T11:25:24.549-03:00"
52     ↪ ^^xsd:dateTime ;
53 19   prov:startedAtTime
54     ↪ "2016-06-16T11:25:12.838-03:00"
55     ↪ ^^xsd:dateTime ;
56 20   wfprov:usedInput
57     ↪ <data/5c65c151-0333/ref/61f8795e-e650>;
58 21
59 22 <data/5c65c151-0333/ref/61f8795e-e650>

```

```

23   tavernaprov:content
24     ↪ <data/4e0baa1f-fc0f/input/
25     ↪ structure.pdb> ;
26 24   wfprov:describedByParameter
27     ↪ <workflow/Setup_MD/processor/split/in/
28     ↪ initial_structure> ;
29 25   wfprov:describedByParameter
30     ↪ <workflow/Setup_MD/processor/in/
31     ↪ initial_structure> ;
32 26   prov:wasGeneratedBy
33     ↪ <run/4e0baa1f-fc0f/process/c3a0e8c0-dcb0/>
34     ↪ ;
35 27   rdf:type wfprov:Artifact ;
36 28   rdf:type prov:Entity .
37 29
38 30 <data/e0fa2f25-0755/ref/55269975-380f>
39 31   tavernaprov:content
40     ↪ <data/4e0baa1f-fc0f/output/bg1c.pdb> ;
41 32   wfprov:describedByParameter
42     ↪ <workflow/Setup_MD/processor/psgen/in/
43     ↪ bg1c_pdb> ;
44 33   wfprov:describedByParameter
45     ↪ <workflow/Setup_MD/processor/split/out/
46     ↪ bg1c_pdb> ;
47 34   wfprov:wasOutputFrom
48     ↪ <run/4e0baa1f-fc0f/process/c3a0e8c0-dcb0/>
49     ↪ ;
50 35   prov:wasGeneratedBy
51     ↪ <run/4e0baa1f-fc0f/process/c3a0e8c0-dcb0/>
52     ↪ ;
53 36   rdf:type wfprov:Artifact ;
54 37   rdf:type prov:Entity .
55 38
56 39 <run/4e0baa1f-fc0f/process/c3a0e8c0-dcb0/>
57 40   wfprov:describedByProcess
58     ↪ <workflow/Setup_MD/processor/split/> ;
59 41   wfprov:usedInput
60     ↪ <data/5c65c151-0333/ref/61f8795e-e650> ;
61 42   prov:wasAssociatedWith <#taverna-engine>
62     ↪ ;
63 43   rdf:type wfprov:ProcessRun ;
64 44   prov:endedAtTime
65     ↪ "2017-03-10T08:19:32.405-03:00"
66     ↪ ^^xsd:dateTime ;
67 45   prov:startedAtTime
68     ↪ "2017-03-10T08:19:31.075-03:00"
69     ↪ ^^xsd:dateTime ;
70 46   prov:used
71     ↪ <data/5c65c151-0333/ref/61f8795e-e650> ;
72 47   wfprov:wasPartOfWorkflowRun
73     ↪ <run/e0fa2f25-0755/> .

```

Listing A.5: Excerpt of the WRO manifest.

```

49 1 @base
50     ↪ <https://w3id.org/w2share/wro/md-setup/>.
51 2 @prefix ro: <http://purl.org/wf4ever/ro#> .

```

```
1 3 @prefix ore:
2   ↪ <http://www.openarchives.org/ore/terms/>
3   ↪ .
4 4 @prefix wf4ever:
5   ↪ <http://purl.org/wf4ever/wf4ever#> .
6 6 <files/Setup_MD/script.sh> a ro:Resource,
7   ↪ wf4ever:Script .
8 7 <workflow/executable-workflow.t2flow> a
9   ↪ ro:Resource, wf4ever:Workflow .
10 8 <workflow/refined-workflow.t2flow> a
11   ↪ ro:Resource, wf4ever:Workflow .
12 9 <data/4e0a1f-fc0f/input/structure.pdb> a
13   ↪ ro:Resource, wf4ever:Dataset .
14 10 <data/4e0a1f-fc0f/output/bglc.pdb> a
15   ↪ ro:Resource, wf4ever:Dataset .
16 11
17 12 <> a ro:ResearchObject ;
18 13   ore:aggregates
19 14     <files/Setup_MD/script.sh>,
20 15     <workflow/executable-workflow.t2flow>,
21 16     <workflow/refined-workflow.t2flow>,
22 17     <data/4e0a1f-fc0f/input/structure.pdb>,
23 18     <data/4e0a1f-fc0f/output/bglc.pdb> .
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
```