

Abstract, Link, Publish, Exploit: An End to End Framework for Workflow Sharing

Daniel Garijo¹, Yolanda Gil¹ and Oscar Corcho²

¹Information Sciences Institute and Department of Computer Science
University of Southern California
{dgarijo,gil}@isi.edu

²Facultad de Informática, Universidad Politécnica de Madrid
ocorcho@fi.upm.es

Abstract

Scientific workflows are increasingly used to manage and share scientific computations and methods to analyze data. A variety of systems have been developed that store the workflows executed and make them part of public repositories. However, workflows are published in the idiosyncratic format of the workflow system used for the creation and execution of the workflows. Browsing, linking and using the stored workflows and their results often becomes a challenge for scientists who may only be familiar with one system. In this paper we present an approach for addressing this issue by publishing and exploiting workflows as data on the Web, with a representation that is independent from the workflow system used to create them. In order to achieve our goal, we follow the Linked Data Principles to publish workflow inputs, intermediate results, outputs and codes; and we reuse and extend well established standards like W3C PROV. We illustrate our approach by publishing workflows and consuming them with different tools designed to address common scenarios for workflow exploitation.

1. INTRODUCTION

Scientific workflows define the "set of tasks needed to manage a computational science process" [Deelman et al., 2009]. They have been used successfully in several domains [Ruiz et al., 2014, Dinov et al., 2009, Wolstencroft et al., 2013] in order to represent, execute, re-run, debug, document and reproduce scientific methods. Scientific workflows are important products of scientific research which should be treated as first-class citizens in cyber-infrastructure [Gil et al., 2007].

There has been great interest in the publication of workflows, particularly to enable reproducibility and shareability of scientific experiments. There are a number of frameworks that allow the inclusion of workflows and codes in scientific articles [Leisch 2007, Mesirov 2010, Falcon 2007]. Workflow repositories like myExperiment [De Roure et al., 2009] and CrowdLabs [Mates et al., 2011] provide mechanisms to publish and search workflows. These repositories support the publication of workflows in their original

language specifications. Given the proliferation of workflow systems both domain-independent [Deelman et al., 2004, Taylor 2006, Scheidegger et al., 2008, Wolstencroft et al., 2013, Ludäscher et al., 2006, Mattmann et al., 2006, Gil et al., 2011] and domain specific [Goecks et al., 2010, Dinov et al., 2009, Reich et al., 2006], this is an impediment for reuse.

An important research challenge is to represent workflows and their associated resources (i.e., inputs, outputs, codes and intermediate results) to facilitate their exploration and adoption by others. Here we tackle this challenge by addressing two main problems: 1) publishing and sharing workflows in a way that can be used by both humans and machines at a low cost; and 2) finding ways to exploit the published data to facilitate reuse and understandability of workflows by other domain scientists.

In order to achieve our goals, this paper describes a framework to publish and exploit computational workflows in a manner that is platform independent. The framework also provides the means to help users understand the published workflows and their associated resources. Our approach makes the following contributions:

- A **collection of requirements and use case scenarios** for workflow reuse and exploitation, based on the state of the art, that we use to assess the features and limitations of our approach.
- A **methodology for publishing workflows as open data on the Web**, in a representation that is independent of the platform used to create them. The methodology publishes workflows according to existing standards and makes them and their associated resources (inputs, intermediate resources, outputs, configurations and codes) available as web objects using the Linked Data principles [Heath and Bizer, 2011].
- A suite of **tools for exploiting and consuming workflow data**, helping end users to overcome the workflow language barrier and address workflow understanding. The suite includes tools for workflow template and execution trace visualization, workflow results browsing, automated documentation of workflows and workflow mining for reuse and revision.

We also present an implementation of this framework for publishing workflows, based on standards like OWL, RDF and PROV [Lebo et al., 2013], and illustrate its use for several workflow domains, ranging from text analytics to ecology. For demonstration purposes, we use the WINGS workflow system¹ [Gil et al., 2011], which has an expressive language to represent reusable abstract workflow templates using semantic constraints in OWL and RDF. Our approach can be applied to other workflow systems, and we illustrate this by publishing workflows from the LONI Pipeline [Dinov et al., 2009].

The rest of the paper is organized as follows. Section 2 introduces prior research on workflow publication and exploitation. Section 3 describes the main scenarios we are interested in addressing, as well as the requirements derived from the scenarios and the state of the art. Section 4 introduces the framework for workflow publication, based on an existing methodology for publishing data on the Web. Section 5 shows how we exploit the data published to facilitate workflow understanding and reuse, while Section 6 discusses the conformance of our framework against the requirements defined on Section 3. Finally, we present our conclusions and future lines of work.

¹ <http://www.wings-workflows.edu>

2. RELATED WORK

Prior related work can be grouped under three major topics: workflow publication, workflow provenance capture and workflow exploitation.

2.1 Workflow Publishing and Sharing

The ability to share workflows among scientists is key for proper reuse and adaptation of existing workflows, as well as a means to provide attribution and feedback from other scientists. The simplest way in which a scientific workflow (or any of its associated resources) may be shared is through data repositories and digital archives. In fact, popular data repositories like Zenodo² and Dryad³ are starting to be heavily used by the community to share datasets and results of their experiments, since they can assign them a Digital Object Identifier⁴ (DOI) for referencing them in their work. DOIs are an important feature for sharing workflow resources, as they ensure their persistence and proper credit through citation. However, in these repositories every description of a resource has to be added manually, and the links and relationships between the rest of the components of the workflow are often missing (e.g., an intermediate result would not be linked to the step that produced it or used it unless it is manually added by the user) even when they have their own entries.

Workflow repositories such as myExperiment [De Roure et al., 2009] and CrowdLabs [Mates et al., 2011] can be used for sharing scientific workflows created with different systems. These repositories store workflows in their native languages, that is, without requiring their conversion to a common language. Although they are great resources for sharing workflows in the community, they don't usually include links to their executions and have to be manually documented. In order to address this issue, some repositories have started to adopt Research Objects [Belhajjame et al., 2015] which bundle together all the resources used within a scientific experiment. However, these are currently in early stages of adoption, relying heavily on manual curation by users.

Another approach for publishing workflows is science gateways, which include applications, libraries of components and workflows that are integrated through a portal. Some examples are described in [Filgueira et al., 2014] for volcanology, in [Danovaro et al., 2014] for hydro-meteorology, or in [Dinov et al., 2009] for neuro-image analysis. Unlike workflow repositories, science gateways usually provide the infrastructure to execute some of the shared components through such community portals.

Finally, prior work [Missier et al., 2010, Shaon et al., 2011] has explored the publication of workflows as Linked Data. In this paradigm, the published workflow resources (e.g., inputs, outputs, intermediate results, etc.) are associated with a URI (Uniform Resource Identifier) to be accessible as web objects. An interesting contribution of that work is to illustrate how the workflow inputs and outputs can be linked to other resources in the Linked Data cloud. However, the workflows (or their resources) were not published using any standard, so they are only reusable by the workflow system used to create them. This

² <https://zenodo.org/>

³ <http://datadryad.org/>

⁴ <http://www.doi.org/>

damages interoperability and affects the potential for understanding and reusability by others.

2.2 Provenance of Workflow Executions

Most of the existing workflow systems capture workflow execution provenance records at some degree in order to help scientists debugging their experiments. Some examples include Galaxy [Goecks et al., 2010], Taverna [Wolstencroft et al., 2013], WINGS [Gil et al., 2011] or Vistrails [Chirigati et al., 2013]. Other tools have also been created to help scientists who use scripting languages like the Jupiter Notebooks⁵ and Apache Zeppelin⁶, which aim to resemble a scientist's lab notebook.

After the consolidation of the W3C standard model for representing provenance in the web (W3C PROV) [Lebo et al., 2013], many systems have started to implement it [Dong et al., 2013], including scientific workflow systems. Different models have been developed to adapt PROV to scientific workflow executions⁷ [Missier et al., 2013, Belhajjame et al., 2015, Garijo and Gil., 2011], but at the moment there is not a standard for scientific workflow representation.

Provenance repositories have started to emerge in order to store workflow executions [Cuevas-Vicentín et al., 2014, Belhajjame et al., 2013]. However, the provenance stored in these repositories is difficult to use at the moment without being familiar with the underlying provenance model.

2.3 Exploitation of Workflow Resources

While a growing number of workflow systems capture workflow executions as provenance records, very few make the associated workflow resources (i.e., inputs, outputs, intermediate resources and codes) accessible and browseable by scientists. The most relevant effort in this regard is the PBase provenance repository [Cuevas-Vicentín et al., 2014], which includes predefined types of queries for the users to edit when exploring the repository. Although PBase is a step forward towards addressing the exploitation of workflow resources, users still need to upload their traces and edit their queries using a graph database language.

General workflow repositories like myExperiment [De Roure et al., 2009] and CrowdLabs [Mates et al., 2011] provide the means to enable workflow accessibility, browsing and search. These repositories have often been exploited by workflow mining researchers to perform different operations with a real workflow corpus for purposes ranging to workflow search enhancement [Starlinger et al., 2014] to workflow clustering [Stoyanovich et al., 2010]. However, they do not contain links to the execution traces.

⁵ <https://ipython.org/>

⁶ <https://zeppelin.apache.org/>

⁷ <http://purl.org/provone>

3. REQUIREMENTS FOR WORKFLOW REUSE AND EXPLOITATION

In this section we review, summarize and expand on the existing requirements defined in the literature to make use of the different aspects of scientific workflows and their executions from an exploitation perspective. The existing infrastructure for the publishing and consumption of provenance traces (and their links to their respective workflow specifications) remains limited to date and does not address these requirements.

3.1 Workflow Reuse and Exploitation Scenarios

We consider here three major purposes for exploiting a workflow repository. One purpose is for gathering information about workflows in the repository in order to understand, evaluate or compare existing results with new work. A second purpose is for updating and adapting existing workflows and their associated resources in order to reuse them to do new research. Both cases are related to workflow reuse, either by the original creators of the workflows, by researchers in their lab or by other colleagues in their institutions.

A third purpose for exploiting a workflow repository is automated workflow mining. Analysis of the workflows stored in the repository can lead to identify groups of results or other salient features of workflows. For example, different efforts have focused on workflow similarity browsing [Bergmann and Gil, 2014], clustering [Stoyanovich et al, 2010], searching [Starlinger et al., 2014] and reuse [Garijo et al., 2013].

We illustrate each of these purposes in the scenarios presented below.

3.1.1 Workflow Information Gathering Scenario

Our first scenario is composed by three user stories. Alice is a proteomics researcher who is working on a new version of a workflow she created previously. Alice wants to evaluate if her new results are consistent with those obtained in her original publication. Luckily, back in the day Alice executed the original workflow several times with different parameter values and datasets, selecting the appropriate results for a publication and storing them in a repository. By looking for the original workflow, Alice is able to recover the successful executions she stored, along with the input datasets and parameter values used to produce them. Alice then uses those input datasets and parameter values with her new version of the workflow and compares the outputs to the originals in order to validate her new workflow. She then publishes her new workflow and its executions in the repository, as a new version.

She is also interested in knowing how her research products (i.e., workflows, software, and data) are used by others within the workflow repository. In particular, she developed a nearest-neighbor clustering algorithm and several workflows to analyze proteomics data and visualize the clustering results. Alice wants this information to report to her funders by pointing to examples where the code, datasets or workflows are used and to collect the list of labs and organizations benefiting from her work. By searching for the name of the software and workflows that she has developed, Alice is able to find a list of the workflow executions where the software and workflows are used. Alice notices that some of those

executions failed due to problems with her software, and recovers the names of their creators and other details that will enable her to track and resolve bugs.

Alice now tries to find similar software that others have developed. She searches for workflows that use hierarchical clustering algorithms, which is the general class of algorithms that include nearest neighbor and other clustering methods. Alice finds hierarchical clustering in abstract descriptions of workflows, and through them she finds several implementations of similar software. Therefore, in order to understand each of the implementations, Alice retrieves the executions corresponding to the abstract workflow and studies their intermediate and final results and the parameter values used as input.

This scenario motivates the following requirements:

- Record a workflow execution and the execution of its steps.
- Retrieval of workflows and associated codes based on their name or identifier.
- Given a workflow, retrieval of its executions.
- Associate inputs, intermediate results, codes and outputs of a given workflow execution.
- Document creators/contributors of a workflow and its executions.
- Document versions of a workflow.
- Document categories of workflow steps.

3.1.2 Workflow Updating Scenario

In our second scenario, Bob is trying to run a workflow for deriving the drug-target network that a student of his lab ran some years ago. Unfortunately, the software that the student used for a molecular docking step is no longer maintained and requires a license that the research lab no longer has. Bob searches the repository for the workflow and finds it, but cannot re-run it with the original settings due to the lack of a license. Bob looks up the workflow template associated with the student's execution specifies abstract steps, and he finds other more recent workflows that execute components of the same abstract class. Thanks to those components, he finds similar software that was used in the last month by another colleague, has the same functionality and is open source. Bob checks the constraints of the input data (e.g., its format and characteristics) and results produced by the new software, and sees that an additional following data conversion step will be needed. Bob updates his workflow with the new step and adds the data conversion step.

The requirements derived from this scenario are the following:

- Retrieval of workflows based on a name or identifier.
- Retrieval of the executions of a workflow.
- Record of the abstract steps of a workflow.
- Retrieval of workflows with a given abstract step.
- Retrieval of Implementations of an abstract workflow step in other workflows.
- Record the types and constraints of data consumed and produced by a step in a workflow.

- Record steps that precede or follow a given step or workflow.

3.1.3 Workflow Mining Scenario

In our final scenario, Clarence is a proteomics scientist that aims to test a hypothesis that a certain protein is present in patient samples for a type of cancer. Clarence knows that his colleagues have published a set of execution runs of different workflows for this kind of analysis, and he wants to exploit those results to test his hypothesis. Thanks to the repository, Clarence is able to cluster all the execution runs based on their common types, parameter values and inputs datasets. After this categorization, he is able to analyze the results in an automated way and assign a confidence value to his hypothesis.

The requirements derived from the scenario are the following:

- Retrieval of the inputs and outputs of a given execution of a workflow.
- Shared inputs and parameter values among groups workflow execution traces.
- Record the types of input data used by a workflow.

3.2 Requirements Summary

The requirements of the previous scenarios can be grouped into four main categories:

- **Workflow template specification requirements**, which tackle the representation and metadata of workflow steps and their dependencies, including how they are connected to each other or any restriction that a given input may have. These requirements also include the metadata of the workflow itself, necessary for tracking authors, contributors, version, license, creation date, etc.
- **Workflow execution (provenance) requirements**, which refer to the processes by which the results of the workflow have been obtained. These set of requirements also refer to the description of the execution as a whole (i.e., its inputs, outputs, intermediate results and codes), including its metadata (authors, contributors, license, etc.). In the literature, this set of requirements has been well defined in the provenance challenge series [Moreau et al 2008].
- **Linkage requirements**, which model the relationships among templates and executions. Some approaches that link both of them together have started to emerge in the last few years⁸ [Missier et al., 2013, Belhajjame et al., 2015, Garijo and Gil., 2011], capturing crucial connections between executions and the methods defining them.
- **Semantic requirements**, which address the types and semantic constraints of data and steps being handled by the workflows. This category of requirements also involves the type of method used in the workflows. Although related work introduces the benefits of this type of semantic information for scientific workflows and tackles the problem from a modeling perspective [Missier et al., 2010, Garijo

⁸ <http://purl.org/provone>

Table 1: List of requirements for workflow template and provenance consumption, grouped by their category and pointing to the scenario where they appear. When an item is in brackets, it indicates that can be referred to any combination of resources included in it.

N°	Category	Requirement	Scenario
R1	Template specification	[templates, steps, inputs, outputs] of the repository with a certain name	Information gathering, updating, mining
R2		Steps that precede or go after a certain step	Information gathering, updating
R3		Steps that [use, generate] a certain variable	Information gathering
R4		[inputs, outputs, intermediate variables, parameters, steps] that belong to a template	Information gathering, updating, mining
R5		[documentation, instructions] defined by users on a template or its components	Information gathering, updating
R6		Current [version, creator, license, contributor] of the template	Information gathering
R7		[creation date, modification date] of a template	Information gathering
R8	Workflow execution	[Inputs, codes, parameters intermediate results] involved in the [execution, generation of a result]	Information gathering, updating, mining
R9		Workflow executions that [used, generated] some [code, inputs, outputs, intermediate results, parameter values]	Information gathering, updating
R10		[Outputs, intermediate results] of workflows with a certain [input dataset, parameter value]	Information gathering
R11		Status of the processes in the run (and the run)	Information gathering, mining
R12		Execution runs that share common configuration parameters, inputs and results	Mining
R13		Creator of the workflow execution	Information gathering
R14		[Start date, end date] of the execution run	Information gathering, mining
R15	Linkage	All executions associated with [a template specification, workflow step, variable]	Information gathering, updating, mining
R16		All templates that are associated with a [input dataset, intermediate result, output]	Information gathering
R17	Semantics	Implementations of an abstract [workflow, component]	Information gathering, updating, mining
R18		Components that have the same function as a given code.	Updating
R19		Workflows that [use, generate] an [input, output, code] type with specific semantic constraints	Updating
R20		All the [datasets, code] of a given type	Updating, mining

and Gil 2011], there are not many workflow systems that are capable of handling these set of requirements for workflow templates and executions.

Table 1 summarizes and generalizes the requirements according to these categories and scenario they belong to. Requirements R1 to R14 may be similar to those tackled in other related work, as they are partially addressed by workflow systems when enacting and executing workflows (requirements R1 to R7) [Wolstencroft et al., 2013, Ludäscher et al., 2006, Mattmann et al., 2006] or representing provenance traces (requirements R8 to R14)¹⁰. Requirements R15 to R20 are the ones where we focus our novel contributions on.

By combining some of the requirement categories (specially the linkage and semantics with the rest), we can create further requirements, such as retrieving the executions of a template that uses or generates certain type of data, comparison of results of runs that used different codes for the same template specification, etc. In Table 1 we have tried to simplify by focusing on those requirements that address the main goals covered by the use case scenarios. An additional list of requirements is available online¹¹.

4. AN END TO END FRAMEWORK FOR WORKFLOW PUBLICATION AND EXPLOITATION

In this section we describe how to publish all workflows and workflow-related products (i.e., the workflow template specifications, execution related resources and metadata) in an automated manner. Several options exist for such publication, as described in Section 2, from workflow repositories to digital archives. In order to be able to reference all the resources properly, we have decided to follow the *Linked Data principles* [Heath and Bizer, 2011]. According to those principles, we should use URIs as names (identifiers) for things, use HTTP URIs so that people can look up those names (making those URIs dereferenceable and available in any browser), provide useful information when someone looks up a URI (by showing the resources that are related to the URI) and include links to other URIs, so anyone can discover additional information.

There are three important advantages of publishing workflows and their resources as Linked Data. The first one is to get linked from other applications by pointing to the URIs that we publish, which include both the workflows and the data generated by them. The second advantage is the ability to produce interoperable results within different systems without having to define particular catalog structures and access interfaces. Just by using standard HTTP operations and formats (like JSON) anyone can use the information published in the repository. Finally, the third advantage is the ability to link to available web resources, for instance referring to proteins in the Protein Data Bank by using their published URI

¹⁰ <http://twiki.ipaw.info/bin/view/Challenge>

¹¹ <https://dx.doi.org/10.6084/m9.figshare.3971994.v1>

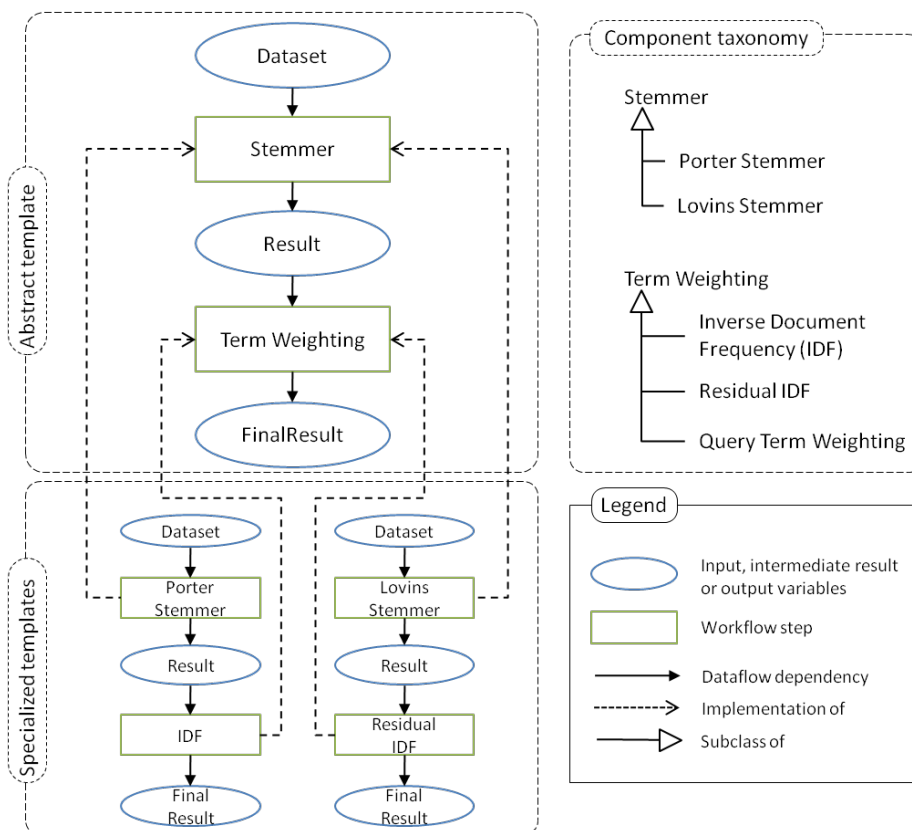


Figure 1: Example of abstract and specialized templates in the WINGS workflow system

This section describes our generic methodology for workflow publication in an end to end scenario. We use WINGS [Gil et al., 2011] as the workflow system of our choice to exemplify our approach. The remainder of the section introduces the key type of abstraction we use to represent workflows, describes the model we have chosen to represent scientific workflows in an infrastructure independent manner (the OPMW model), illustrates the architecture and methodological steps that need to be carried out to publish the workflow template specifications and their executions as Linked Data on the Web and summarizes the main features of WINGS.

4.1 Representing Semantic Abstractions

We represent workflows with skeletal plan abstractions [Friedland and Iwasaki 1985], where every step in the abstract plan always corresponds to a class of steps in the concrete plan. This type of abstraction allows a simple generalization of the steps in the workflow. An example can be seen in Figure 1, showing an *abstract* template on the top and two *specialized* (concrete) templates on the bottom. By looking at the taxonomy of workflow steps shown in the top right of the figure, we can see that both of the specialized templates shown are implementations of the same abstract method (with the Stemmer and Term Weighting steps).

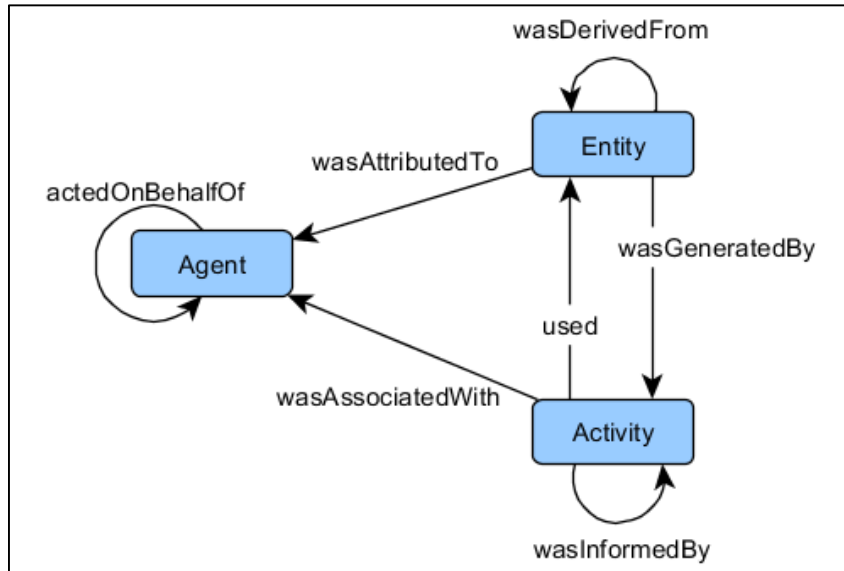


Figure 2: Overview of PROV

This way, an abstract workflow template may have different levels of abstraction, depending on the generality aimed at by the designer on the workflow. This approach also supports mixing specific step implementations within an abstract template (e.g., if a particular step of the workflow will always use the same implementation).

4.2 Workflow Template and Execution Representation: OPMW

As we mentioned in Section 2, W3C PROV is a standard for representing provenance on the Web. PROV uses three main concepts to represent provenance: *Entities*, i.e., the things we want to describe; *Activities*, i.e., the processes that consumed or produced entities; and *Agents*, who are the entities responsible for carrying out the activities. As shown in Figure 2, these concepts capture the provenance of an entity by using seven main relationships: *used* (an activity used some entity), *wasAssociatedWith* (an agent participated in some activity), *wasGeneratedBy* (an activity generated an entity), *wasDerivedFrom* (an entity was derived from another entity), *wasAttributedTo* (an entity was attributed to an agent), *actedOnBehalfOf* (an agent acted on behalf of another agent) and *wasInformedBy* (an activity uses results from another activity).

PROV was designed to be generic and domain independent, and needs to be extended to address the requirements to represent workflow templates and executions. In order to achieve this goal, we have developed the OPMW model¹² [Garijo and Gil, 2011]. Since PROV only has one term to refer to the plan associated with the execution of an activity (called *Plan*), OPMW adopts the P-Plan¹³ model [Garijo and Gil 2012], a model designed to capture the steps and variables of scientific processes, for addressing the template representation requirements. In addition, OPMW also extends the Open Provenance Model

¹² <http://www.opmw.org/ontology/>

¹³ <http://purl.org/net/p-plan>

(OPM) [Moreau et al., 2011], a legacy provenance model developed by the workflow community that was used as reference to create PROV.

OPMW supports the representations of workflows at a fine granularity with a lot of details pertaining to workflows that are not covered in those more generic languages. OPMW also allows the representation of links between a workflow template and a workflow execution that resulted from it. Finally, OPMW also supports the representation of roles and attribution metadata about a workflow.

4.2.1 OPMW Overview

In OPMW, a workflow template specification is represented as an *opmw:WorkflowTemplate*. Workflow templates are formed by *opmw:WorkflowTemplateProcesses*, which use or produce *opmw:WorkflowTemplateArtifacts* hence capturing the dataflow dependencies.

On the execution side, each *opmw:WorkflowExecutionProcess* represents the execution of a workflow template process, and is bound to it via the *opmw:correspondsToTemplateProcess* relationship. Similarly, each *opmw:WorkflowExecutionArtifact* that is used or generated by a workflow execution process is linked to its corresponding workflow template artifact with the *opmw:correspondsToTemplateArtifact* relationship. Finally, the *opmw:WorkflowExecutionAccount* containing all the provenance statements of the execution is linked to the workflow template that contains all the assertions of the template with the *opmw:correspondsToTemplate* relationship.

Figure 3 shows an example of the OPMW vocabulary extending OPM, PROV and P-Plan. Each vocabulary concept and relationship is represented with a prefix to help understand its source (OPMW, P-Plan, PROV or OPM (with opmo and opmv)). In the figure, an abstract workflow template with one sorting step, an input and an output (on the top right of the figure, represented using P-Plan) is linked to its provenance trace on the bottom right of the figure (depicted with PROV and OPM). Each activity and artifact of the execution is linked to its respective step and variable on the template. Additional metadata of the variables (e.g., constraints), steps (e.g., conditions for execution), activities (e.g., used code), artifacts (e.g., size, encoding), account (e.g., status) and template (e.g., associated dataflow graph) is modeled with OPMW, but has been omitted from the figure for simplicity.

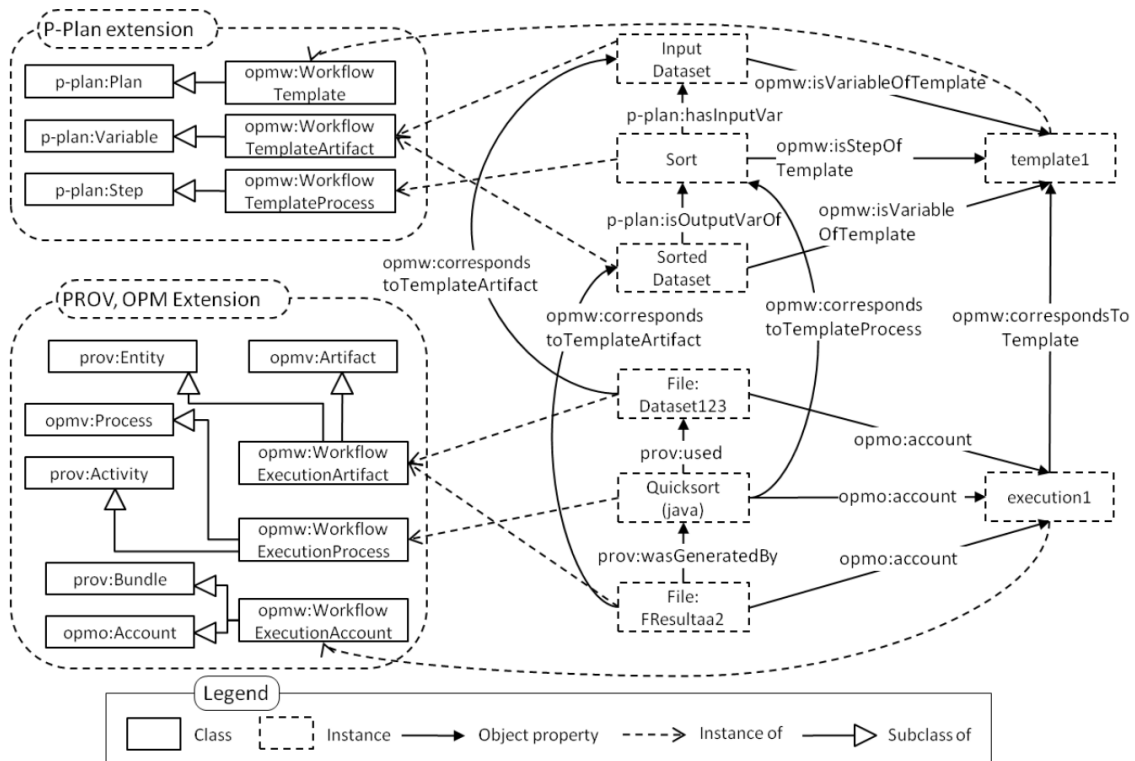


Figure 3: OPMW example of an execution (bottom right of the figure) and the links to its corresponding workflow template (top right of the figure). The extensions to P-Plan, OPM and PROV can be seen in the left part of the figure.

Attribution is crucial for scientists who create and publish workflows, as well as for those who provide the data or workflow execution infrastructure. Hence, OPMW reuses terms from existing vocabularies to represent them. For example, the Dublin Core (DC) Metadata Vocabulary is reused to represent the author, contributor, rights and license of an input dataset, a code used in a workflow or the workflow itself. OPMW also defines additional terms for referring to the start and end of the whole execution of the workflow, the size of the produced files, the status of the final execution, the tool used to design the workflow, the tool used to execute the workflow, etc.

4.3 Publishing Workflows as Data on the Web

There are methodologies for Linked Data generation and publication [Heath and Bizer, 2011]. However, there are no methodologies for publishing workflows and their constituents online. Therefore, we have developed a methodology that adapts an existing one used in the publication of government data [Villazon-Terrazas et al., 2012] and other domains like energy consumption [Radulovic et al., 2015] to achieve our purpose.

The methodology consists of five main steps:

- 1) **Specification**, where the of data sources to use are identified and a URI naming convention is designed (the license for the resulting dataset is often agreed during this step as well);

- 2) **Modeling**, where users decide which vocabularies should be used to represent the data properly according to the existing requirements and scenarios;
- 3) **Generation**, i.e., the process of transforming the data from their heterogeneous formats to a structured format, typically RDF¹⁴, with the help of existing tools; cleaning the data and linking it with other existent sources
- 4) **Publication**, where the resulting dataset and its metadata is made available by using an online accessible triple store; and
- 5) **Exploitation**, where the benefits of the dataset is made clear through applications or queries that consume it.

An overview of the architecture for our implementation of the methodology can be seen in Figure 4. For the specification step, we decided that all the URIs generated by our system would become *cool URIs*¹⁵, following the W3C recommendations. This means that they are produced under a domain under our control, they are unique, and they are not going to change. Each URI identifies a different resource that can be individually accessed and dereferenced with content negotiation (i.e., the same URI can handle requests for users (returning HTML) and machines (returning an RDF serialization)). The URIs adopt the following naming scheme:

Base URI = <http://www.opmw.org/>

Ontology URI = <http://www.opmw.org/ontology/>

Assertion URI = <http://www.opmw.org/export/resource/ClassName/instanceID>

As for the license, we choose the Creative Commons Attribution-Share Alike 3.0¹⁶, which allows anyone to reuse the published data if proper attribution is provided. However, the license might be changed if the authors prefer another one.

For the modeling step, we chose the OPMW vocabulary. Other approaches like the Research Object Model [Belhajjame et al., 2015] or the PROV-One vocabulary¹⁷ may be considered when addressing the modeling step.

For the generation step, the workflow template specifications and execution traces are converted automatically to OPMW with a transformation script¹⁸. The script takes as input a workflow template or an execution of the target workflow system, and produces an RDF file using the naming convention established in the first step. Camel case notation is used for composing the identifiers of classes and instances, and an MD5 encoding is used to generate a unique identifier for each resource when necessary (e.g., a used file, an execution step, etc.).

¹⁴ <https://www.w3.org/RDF/>

¹⁵ <http://www.w3.org/TR/cooluris/>

¹⁶ <http://creativecommons.org/licenses/by-sa/3.0>

¹⁷ <http://purl.org/provone>

¹⁸ <https://github.com/dgarijo/WingsProvenanceExport>

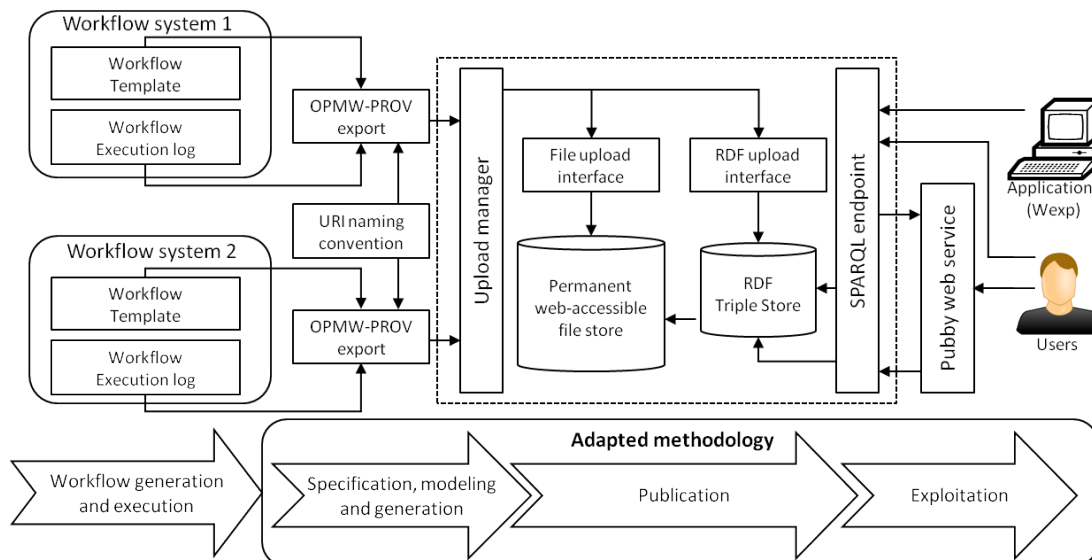


Figure 4: Overview of our architecture for publishing workflow as Linked Data

The **publication step** takes as input the RDF files produced by the generation step. As shown in Figure 4, an upload manager loads the files into a triple store and makes them available through a public endpoint (i.e., an access point for both human users and machines). We have selected Openlink Virtuoso¹⁹ as our triple store because of its robustness, support from the community and mechanisms to create a public access. Whenever an execution run is published, all the referenced resources (inputs, intermediate results and outputs of the workflows) are versioned and saved in an additional permanent file store. If the file size is too big (size can be configured), then the pointers to the original location are maintained. The upload interface makes sure the consistency with the links used in the execution traces in the triple store, so they can be accessed appropriately. The file store is available in our local servers (<http://www.opmw.org>). A Linked Data Frontend is set up using Pubby²⁰ for facilitating the access to each individual resource of the endpoint. The public endpoint is available online²¹, along with a set of sample queries to retrieve basic data from workflows and demonstrate its main functionality²².

The last step of methodology, the exploitation phase, is described in Section 5.

4.4 Designing and Executing Workflows: WINGS

WINGS is a workflow system that uses semantic representations to describe the constraints of the data and computational steps in the workflow [Gil et al., 2011, Gil et al., 2011b]. WINGS can reason about these constraints, propagating them through the workflow structure and using them to validate workflow templates. It has been used in different domains, ranging from life sciences to multimedia analysis and geosciences.

¹⁹ <http://virtuoso.openlinksw.com/>

²⁰ <http://wifo5-03.informatik.uni-mannheim.de/pubby/>

²¹ www.opmw.org/sparql

²² <http://www.opmw.org/node/6>

WINGS provides web-based access and can run workflows locally, or submit them to execution environments such as Pegasus/Condor [Deelman et al., 2005] or Apache OODT [Mattmann et al., 2006], so as to handle large-scale distributed data and computations, optimization and data movement.

There are two main reasons why we have selected WINGS to exemplify our methodology for publishing workflows and their resources on the Web. The first one is that WINGS separates types of data and types of components when designing the workflow. This results in a data catalog and a component catalog organized in a taxonomy, which we can exploit for addressing the semantic requirements depicted in our scenarios. The second reason is that WINGS allows creating abstract workflow templates [Gil et al., 2011], similar to those presented in Section 4.1, which can then be instanced differently on each execution.

5. WORKFLOW DATA EXPLOITATION

Once the workflow related data are available as web resources, they can each be independently accessed with their metadata, or through queries to the endpoint to retrieve more complex relationships. The data consumption can be done both from a human perspective and a machine point of view, since all the contents are derreferenceable and offer serializations in multiple formats on demand.

However, for "true users" [Cohen and Leser 2011] like Alice, Bob or Clarence this is often not enough, as they lack the skills necessary to be familiar with query languages like SPARQL, SQL or similar and they do not know the models or schemas the workflow data is represented in.

In this section we describe our approach for addressing this issue. We present a suite of tools (developed by ourselves or third parties) that we use for exploiting the contents of a repository of workflow templates and their execution traces published according to the principles discussed in Section 4.

5.1 Workflow Information Gathering and Adaptation

We help end users browse, inspect, documentation and visualize the contents of the repository in an easy manner, for that does not require that they issue any queries against the workflow repository.

5.1.1 Workflow Browsing

The Workflow Explorer (WExp)²³ allows navigating over different workflow templates, their metadata and their workflow execution results. WExp loads dynamically the workflow information stored in a repository, and allows the user to search it on demand.

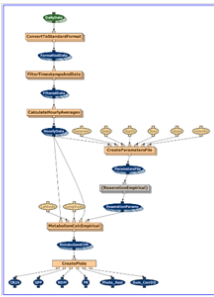
²³ <https://github.com/dgarijo/WorkflowExplorer>

Workflow browser

Type the name or the initial letters of the template you want to browse and the system will autocomplete it. For example, type "a" in the box below:

AquaFlow_EM

Template metadata [-]	
Template modified on:	undefined
Label:	AquaFlow_EM
Diagram URI:	AquaFlow_EM.owl.png
Created in workflow system:	http://wings.isi.edu
Native system template URI:	AquaFlow_EM.owl
Version number:	2
Template contributor:	http://www.opmw.org/export/resource/Agent/WATER



Data Variables [-]	Parameter Variables [-]
<ol style="list-style-type: none"> FormattedData (more) ParametersFile (more) FilteredData (more) MetabolismEDM (more) DailyData (more) PR (more) ReaerationParams (more) HourlyData (more) NDM (more) GPP (more) CR24 (more) Sum_CorrDO (more) Photo_Rest (more) 	<ol style="list-style-type: none"> date (more) depth (more) velocity (more) flow (more) slope (more) Latitude (more) barpress (more) Longitude (more)

Available executions [-]	Processes [-]
<ol style="list-style-type: none"> ACCOUNT1383026688239 ACCOUNT1383027378682 	<ol style="list-style-type: none"> ConvertToStandardFormat (more) MetabolismCalcEmpirical (more) CreateParametersFile (more) ReaerationEmpirical (more) FilterTimestampsAndData (more) CalculateHourlyAverages (more) CreatePlots (more)

Figure 5: Snapshot of WExp showing the information of a workflow.

The entry point for the application is a workflow template name. Whenever a user like Alice or Bob selects the one they are looking for, a series of lists with information appear on the screen as shown in Figure 5. Each list shows resources related to the template, grouped by their common type and retrieving the data asynchronously from the server. First, the metadata list of the template is shown. It includes the creators and contributors (such as the ones Alice was searching to include in her report), the template license, version number, and the workflow system used to create it, as well as a picture of its workflow template specification (if available).

The rest of the lists show the data variables of the workflow template (i.e., inputs, outputs and intermediate variables) and their constraints, the steps of the workflow (so users like Bob can look for their configuration and results on the execution traces) and the available execution traces (which Alice wanted to explore for testing the new version of her workflow). Each resource is a link that can be resolved in the browser for more information, making it easy to know if an execution run has been successful or not, find the available implementations of a workflow step or how a variable in the workflow template has been instanced in the available executions. A live demo of WExp can be found online²⁵.

²⁵ <http://purl.org/net/wexp>

The screenshot displays the 'AQUAFLOW EM' workflow page. On the left is a navigation sidebar with 'Organic Data Publishing' logo and links like 'Main page', 'Recent changes', and 'Tools'. The main content area is titled 'AQUAFLOW EM' and includes sections for 'Workflow' (listing 'AQUAFLOW EM'), 'Processes' (listing 'CREATEPARAMETERSFILE', 'CALCULATEHOURLYAVERAGES', 'FILTERTIMESTAMPSANDDATA', 'CONVERTTOSTANDARDFORMAT', 'METABOLISM CALCEMPIRICAL', 'REAERATIONEMPIRICAL', 'CREATEPLOTS'), 'Workflow Executions' (listing 'AF_EM_Execution_2_March_2012_to_8_March_2012'), 'Contributor' (listing 'WATER'), 'Workflow Created In' (listing 'wings.isi.edu'), and 'Template File' (listing 'AquaFlow_EM.ow').

The top right section, titled 'AF EM Execution 2 March 2012 to 8 March 2012', shows 'Executed Workflow' (ACCOUNT1383027378682), 'Input Data' (AQUAFLOW_EDM_DAILYDATA, DAILYDATA, AQUAFLOW_NTM_DAILYDATA), and 'Generated Data' (NEM, METABOLISM, AQUAFLOW_EDM_HOURLYDATA, HOURLYDATA, AQUAFLOW_NTM_HOURLYDATA, AQUAFLOW_EDM_FILTEREDDATA, FILTEREDDATA, AQUAFLOW_NTM_FILTEREDDATA, AQUAFLOW_EDM_FORMATTEDDATA, FORMATTEDDATA, AQUAFLOW_NTM_FORMATTEDDATA, REAERATIONPARAMS, AQUAFLOW_EDM_PARAMETERSFILE, PARAMETERSFILE, AQUAFLOW_NTM_PARAMETERSFILE, SUM_CORRDO, CR24, PHOTO_REST, PR, GPP). A black arrow points from 'CREATEPLOTS' in the workflow to 'GPP' in the generated data.

The bottom right section, titled '995dfbd9728f3fd06979ecf14a3e2ccd', shows 'Properties' for the 'GPP' output:

Bapress	760
Depth	1.04
Flow	1,581,684
ForDate	3 March 2010
ForSite	SMN
HasSize	6,163
SiteLatitude	37.347
SiteLongitude	-120.976
Slope	1.0e-4
Velocity	0.662
WasGeneratedBy	http://www.opmw.org/portal/resource/WorkflowExecutionProcess/CREATEPLOTS/1383027378682

The 'Credits' section at the bottom right lists 'Users who have contributed to this Page' as 'Admin (15 Edits)'.

Figure 6: Snapshot of the Organic data publishing wiki, showing a page of a workflow template (left), its execution (top) and metadata of one of the used values (bottom right)

5.1.2 Workflow Documentation

The Organic Data Science Wiki²⁷ is an extension of semantic wikis designed to develop meta-workflows that result in many workflow explorations and runs. A snapshot of the interface is shown in Figure 6, depicting a workflow template on the left, one of its executions on the top and the detail for one of the generated outputs on the bottom. By using this wiki, end users can import workflow templates and workflow executions into this framework to generate persistent documentation pages that link to data and algorithm descriptions in the wiki. End users can augment this documentation as they relate to the overall meta-workflow. This may help not only improve the quality of the workflows and results they use for their experiments, but also foster collaboration with other scientists in other domains. The pages look similar to what WExp presents, as they include inputs and outputs, metadata of the workflow, steps and codes used in the executions, etc. The advantage of this approach is that the initial page for the workflow documentation is populated automatically from the contents of the repository, making it easy to navigate through the workflow and explore its main resources and executions. While WExp is for searching and retrieving particular types of data, the Organic Data Science Wiki focus on high quality descriptions and automated documentation of workflows, their execution traces and results.

²⁷ http://www.organicdatacuration.org/sandbox/index.php/Main_Page



Figure 7: PROV-O-Viz sample visualization diagram of a workflow execution

Finally, the components used for the workflow execution can be further browsed and documented in the OntoSoft portal³¹ [Gil et al., 2015]. This is especially useful for users who want to find and replace old software with other compatible versions equal in functionality. The OntoSoft portal is a platform created to provide a community repository for software, allowing contributors to specify structured metadata about software, assisting users to make their software more reusable (e.g., by suggesting which information is missing to be able to rerun it) and maintaining compatibility with existing software catalogues such as Github or BitBucket. The OntoSoft ontology [Gil et al., 2015]³² is used by the portal to provide a vocabulary for allowing users to publish and describe their software in terms of accessibility, execution, citation and understanding.

5.1.3 Workflow Execution Trace Visualization

We use the PROV-O-Viz tool [Hoekstra and Groth 2014] for visualizing workflow executions stored in the repository, so users like Alice can get a quick overview of an execution trace. PROV-O-Viz creates a Sankey³³ diagram out of the execution trace of the workflow, which helps provide a general insight into how the inputs influence the final results. A screenshot is shown in Figure 7, showing a sample of a simple workflow execution. PROV-O-Viz was designed to work in a generic domain with the W3C PROV standard, but we consider a useful asset for creating overviews of workflow traces. Thanks to the OPMW compatibility with PROV, we are able to benefit from this and other applications that may be developed in the future.

5.2 Mining Workflow Data

The efforts shown in the previous section focus on presenting and facilitating the access to data to users so they can consume it without having to perform complex queries to a repository. Next, we detail our own efforts to consume the data for exploiting its value.

5.2.1 Meta-workflow Analysis

An important aspect of publishing workflows is being able to consume them afterwards to review and use their findings for other experiments. We do so in the DISK system [Gil et al 2016],³⁴ introducing a novel approach to automate the hypothesize-test-evaluate

³¹ <http://www.ontosoft.org/portal/>

³² <http://ontosoft.org/software#>

³³ <http://bost.ocks.org/mike/sankey/>

³⁴ <http://disk-project.org/>

discovery cycle with an intelligent system that a scientist can task to test hypotheses of interest in a data repository.

The DISK approach captures three types of data analytics knowledge: 1) common data analytic methods represented as semantic workflows; 2) meta-analysis methods that aggregate those results, represented as meta-workflows; and 3) data analysis strategies that specify for a type of hypothesis what data and methods to use, represented as lines of inquiry.

Given a hypothesis specified by a scientist, appropriate lines of inquiry are triggered, which lead to retrieving relevant datasets, running relevant workflows on that data, and finally running meta-workflows on these and previous workflow results (retrieved from the repository). The scientist is then presented with a level of confidence on the initial hypothesis (or a revised hypothesis) based on the data and methods applied.

5.2.2 Finding Commonly Occurring Fragments in Workflows

We aim to facilitate workflow reuse and understanding by mining the templates of our repository to find commonly occurring fragments of workflows. Our approach, FragFlow [Garijo et al., 2013], collects all the workflow templates from an existing repository and performs a series of analyses by using state of the art frequent sub-graph mining techniques. As a result, we obtain a series of *workflow fragments*, which relate workflows to each other and indicate the parts of the repository that are more likely to be reused. FragFlow also exploits the abstract methods available in the repository, verifying whether two templates might share common fragments at a higher level of abstraction. When applied to multi-disciplinary workflow repositories, FragFlow may be used to discover common methods used by different users, groups or even disciplines.

6. RESULTS AND DISCUSSION

Having introduced our suite of tools, in this section we aim to assess (1) the feasibility of our methodology for publishing workflows in the web in an automated manner, (2) whether our suite of tools fully addresses the requirements identified in Section 3 for workflow exploitation and (3) additional exploitation opportunities that we have not addressed in this work.

6.1 Methodological Feasibility

Following the steps introduced in Section 4, we have published on the public endpoint 73 workflow templates and 106 workflow executions traces of the WINGS workflow system. In addition, in order to test the feasibility of our approach with other workflow systems, we have also published 91 workflows from the LONI Pipeline public library.³⁵ The workflows belong to different domains, ranging from text analytics to pharmacology or water quality analysis. The size of each workflow is also variable. For example, we have workflows with only a few steps, while others have executions with more than 170. Each template includes

³⁵ <http://pipeline.loni.usc.edu/explore/library-navigator/>

all its original metadata and links to its corresponding executions (if any). For each execution, all the metadata and inputs, results, intermediate outputs and used codes have been made available as well. This includes files ranging from a few bytes to almost 20 GB. All the workflows and executions, as well as the tools presented in previous sections, have been made part of the *WEST workflow ecosystem* [Garijo et al., 2014], which integrates them through OPMW and P-Plan to interoperate at different levels of abstraction.

From the modeling perspective, we have tested the feasibility of our approach by answering each of the requirements listed in Table 1 with the OPMW model and examples from the repository. Each requirement has been assigned a SPARQL query that addresses it. The resultant table from this effort is available online.³⁶ Since each of the requirements can be answered successfully, we conclude that our model is feasible to represent the workflow templates and executions we want to exploit.

Regarding the validation of the resources required to support publication according to the Linked Data principles, we have used the Vapour system,³⁷ which launches a set of tests retrieving the exposed data in different formats typically consumed by humans (e.g., html) and machines (e.g., RDF/XML). In addition, we have performed a set of unitary tests to check any inconsistencies on the data.³⁸

Given that we are able to answer the requirements for workflow template and execution exploitation, we have published successfully a corpus of workflows and executions and we have validated that they are accessible under the Linked Data principles, we can conclude that our proposed methodology is sufficient for publishing workflows as data on the web.

That said, keeping track and storing all the data resultant from an experiment may cause scalability issues in certain scenarios. For example, external repositories may be used as an input of a workflow. In other cases, thousands of execution runs may have to be explored before achieving the desired results. In our approach, the total amount of triples required by the system to represent the 270 workflow templates and executions in the repository is around 136,000. Triple stores are nowadays capable of handling billions of triples, which would ensure the storage of tens of thousands of workflow executions and templates in a single instance. Regarding the accessibility to external repositories with big data, we don't consider a big issue not storing the whole dataset each time an experiment is run. Instead, the author can provide enough metadata by, for example, linking to the source repository and to the queries used to retrieve the workflow data.

6.2 Workflow Exploitation: Conformance to Requirements

In this section we analyze whether each of the requirements defined in Section 3 are addressed by the tools we propose to facilitate the exploitation of workflow data. The first category is the template specification requirements. Thanks to the Workflow Explorer and the Organic Data Science Wiki, template metadata like documentation (R5), version, creator

³⁶ <https://dx.doi.org/10.6084/m9.figshare.3971994.v1>

³⁷ <http://linkeddata.uriburner.com:8000/vapour>

³⁸ <https://github.com/dgarijo/OPMWValidator>

(R6) and creation dates (R7) are shown to the user when accessing a workflow by its name (R1). Similarly, the inputs, outputs and other resources (R4) are also displayed when a workflow template is accessed in any of the applications. The steps that use or generate a certain variable within the workflow template (R3) can be solved by resolving the URI of the desired template in the browser. Finally, the list of steps that precede or go after a certain step (R2) is not directly shown in the applications, but can be inferred easily from the list of processes and the available visualization of the workflow.

The workflow execution requirements are addressed in a similar way. The metadata of the execution e.g., status (R11), creator (R13) and start and end dates (R14) are shown by default whenever accessing an execution of the Organic Data Science Wiki, along with the codes, inputs and intermediate values that are involved in the generation of a result (R8). The workflow executions that used or generated a certain dataset (R9) and the outputs of workflows with a certain input (R10) can be found out by searching in the wiki for the result or by simply resolving the result itself to obtain an HTML description of its metadata. However, the current support for finding groups of execution runs that share parameter values or inputs (R12) is limited (not directly shown on our current suite of tools), and part of our future work.

The remaining two categories of requirements are addressed by WExp, the Wiki and our mining applications. For each template searched in the explorer, a link to each execution is provided (R15), and whenever a resource URI is resolved against the server, the templates that instantiates it are specified with the *correspondsToTemplateArtifact* relationship (R16). The implementations of an abstract component (R17) can be seen by dereferencing the URI of that component, by looking for the workflow in WExp or by running FragFlow in the repository to find abstract workflow fragments. Abstract workflow fragments link together the different parts of workflows which share the same functionality, being able to know which components share the same function at a higher level of abstraction (R18). Finally, by looking for a certain semantic type on the Wiki, one can retrieve files associated with it (R20), which can be further explored to retrieve the workflow they belong to (R19).

In summary, 18 of 20 the requirements are fully supported by our suite of tools, and when they are not supported directly, the applications provide the information necessary to address them.

6.3 Additional Exploitation Opportunities

There are two main additional aspects that have not been addressed in our work, but that may benefit from our current approach. The first one is the automatic re-execution of the published workflows. This is facilitated because the provenance traces have all the required information and code used in the workflow, while the workflow templates represent the order and data dependencies needed to link the different steps. In fact, we have successfully used OPMW and P-Plan to execute published workflows in different execution engines [Gil 2013]. However, executing workflows in a different environment would require us to extend our current approach: important metadata would need to be captured, such as the software dependencies that have to be installed for each of the workflow steps or the execution

environment requirements (memory, OS, CPU, etc.). We believe that by incorporating models that capture the missing metadata (e.g., WICUS [Santana-Pérez and Pérez-Hernandez, 2015]) or by using virtual execution environments (e.g., Docker containers³⁹) we may be able to address most of the challenges associated with workflow re-execution.

The second main aspect that would benefit from our work is the automatic generation of explanations for different workflow results. We believe that despite all the tools presented here for workflow visualization and browsing, there is still a need to summarize, compare and explain the results of the workflow in a way that is human readable, depending on the users' needs (e.g., an abstract method for understanding the main steps, the specific software used, etc.). Our approach creates the knowledge base needed to support these kinds of explanations and interactions between different levels of abstraction.

7. CONCLUSIONS AND FUTURE WORK

An important movement in science is open science, increasingly facilitating access to scientific products by anyone. By analogy with open publications and open data, the open availability of workflows would enable greater reuse of workflows as scientific artifacts. This often involves the publication of the workflow associated resources, in order to provide a provenance trace with evidence of how a particular result was obtained.

In this paper we have described our approach to publish abstract workflow templates and their respective workflow execution traces as data on the web. We believe that publishing workflows and execution traces without providing the means to end users to consume these resources may be a barrier for their reuse. Our contributions to address this issue include a collection of requirements and use cases for consuming and exploiting workflow resources (derived from our own experience and the state of the art), a methodology for publishing workflow as web objects and a suite of tools that can consume the published data in order to facilitate its understanding to users.

The benefits of our approach promote the ability to share and consume each of the published resources by using standard HTTP operations, allow the definition of an evaluation framework made from the requirements collected from the use case scenarios and enable the interoperability between the applications exploiting the data and third party tools. In addition, thanks to our representation based on standards, the published workflows and their traces are agnostic about the workflow system used for their design.

We have tested the feasibility of our methodology by publishing a corpus of workflows from two workflow systems and we have evaluated our approach against the requirements driven by the common scenarios. As a result, we have verified that the majority of the requirements are satisfied by our current suite of applications. Those requirements only partially addressed (e.g., comparison among workflow execution traces) are part of our ongoing work.

³⁹ <https://www.docker.com/what-docker>

Regarding future work, our current approach has focused on the publication of all the resources involved in the workflow as Linked Data on the Web. However, some of these resources may contain sensitive information that should not be published along with the rest of the experiment. Although the endpoint and server storing the resources can be modified to limit the access to certain published files, we are currently studying the alignment of our approach with a conditional access Linked Data framework⁴⁰ that allows users accessing on the data based on their roles and privileges. On a related point, we also aim to explore how our publishing framework could be used in the early stages of workflow development (i.e., sharing the preliminary contents of the research only with the working group).

Another area of related work relies on the interoperability of the represented data. As we have shown in our previous work [Garijo et al., 2014] the OPMW model we have used to represent the workflow data can be mapped to other existing approaches. Since there are emerging efforts that publish workflow data according to these other approaches [Cuevas-Vicentín et al., 2014] we aim to be able to exploit these data as well. In addition, we plan to explore the possibility of having orchestrating services that help the communication and set up of the workflow ecosystem comprising all the tools consuming the workflow data (e.g., those used for DevOps in software engineering projects).

Finally, another challenge to be addressed is regarding the persistence of the published resources. At the moment we assign them URIs that aren't expected to be changed, but we would like to associate each of the resources with a DOI so as to be able to cite them in publications.

ACKNOWLEDGEMENTS

The authors would like to thank Daniel Vila for his feedback on this work. We also gratefully acknowledge support from the Defense Advanced Research Projects Agency through the SIMPLEX program with award W911NF-15-1-0555, and the US National Science Foundation with award ICER-1440323.

REFERENCES

- [Belhajjame et al., 2013]: Khalid Belhajjame, Jun Zhao, Daniel Garijo, Aleix Garrido, Stian Soiland-Reyes, Pinar Alper and Oscar Corcho. A Workflow PROV-Corpus Based on Taverna and WINGS. Proceedings of the Joint EDBT/ICDT 2013 Workshops, pages 331-332. Genova, Italy 2013.
- [Belhajjame et al., 2015]: Belhajjame, K., Zhao, J., Garijo, D., Gamble, M., Hettne, K., Palma, R., Mina, E., Corcho, O., Gómez-Pérez, J. M., Bechhofer, S., Klyne, G., and Goble, C. Using a suite of ontologies for preserving workflow-centric Research Objects. Web Semantics: Science, Services and Agents on the World Wide Web. 2015.
- [Bergmann and Gil, 2014]: Bergmann, R. and Gil, Y. (2014). Similarity assessment and efficient retrieval of semantic workflows. Information Systems, 40:115-127. 2014.
- [Chirigati et al., 2013]: Chirigati, F., Freire, J., Koop, D., Silva, C. VisTrails provenance traces for benchmarking. in Proceedings of the Joint SDBT/ICDT 2013 Workshops, 323-324, 2013.

⁴⁰ <http://conditional.linkeddata.es/>

- [Cohen and Leser, 2011]: Cohen, S. and Leser, U. Search, adapt and reuse: the future of scientific workflows. *ACM SIGMOD Record*, 40:2, 6-16- 2011.
- [Cuevas-Vicenttín et al., 2014]: Víctor Cuevas-Vicenttín, Parisa Kianmajd, Bertram Ludäscher, Paolo Missier, Fernando Chirigati, Yaxing Wei, David Koop, Saumen Dey. The PBase Scientific Workflow Provenance Repository. *International Journal of Digital Curation* 9(2): 28-38, 2014.
- [Danovaro et al., 2014]: Danovaro, E., Roverelli, L., Zereik, G., Galizia, A., DAgostino, D., Paschina, G., Quarati, A., Clematis, A., Delogu, F., Fiori, E., Parodi, A., Straube, C., Felde, N., Harpham, Q., Jagers, B., Garrote, L., Dekic, L., Ivkovic, M., Caumont, O., and Richard, E. (2014). Setting up an hydro-meteo experiment in minutes: The DRIHM e-infrastructure for HM research. In *e-Science (e-Science)*, 2014 IEEE 10th International Conference on, volume 1, pages 47-54
- [Deelman et al., 2004]: Deelman, E., Blythe, J., Gil, Y., Kesselman, C., Mehta, G., Patil, S., Su, M. H., Vahi, K., and Livny, M. "Pegasus: Mapping scientific workflows onto the grid". In Dikaiakos, M., editor, *Grid Computing*, volume 3165 of *Lecture Notes in Computer Science*, pages 11- 20. Springer Berlin / Heidelberg. 2004.
- [Deelman et al., 2009]: Deelman, E., Gannon, D., Shields, M., and Taylor, I. Workflows and e-science: An overview of workflow system features and capabilities. *Future Generation Computer Systems*, 25(5):528-540. 2009.
- [De Roure et al., 2009]: De Roure, D; Goble, C.; Stevens, R. "The design and realizations of the myExperiment Virtual Research Environment for social sharing of workflows". *Future Generation Computer Systems*, 25 (561-567), 2009.
- [Dinov et al., 2009]: Dinov, I. D., Horn, J. D. V., Lozev, K. M., Magsipoc, R., Petrosyan, P., Liu, Z., MacKenzie-Graham, A., Eggert, P., Parker, D. S., and Toga, A. W. Efficient, distributed and interactive neuroimaging data analysis using the LONI Pipeline. In *Frontiers in Neuroinformatics*, volume 3. 2009.
- [Dong et al., 2013]: Dong-Huynh, T.; Groth, P.; Zednik, S. PROV Implementation Report. W3C Working Group Note. WWW Consortium. 30 April 2013.
- [Falcon 2007]: Falcon, S. "Caching code chunks in dynamic documents: The weaver package." *Computational Statistics*, (24) 2, 2007.
- [Filgueira et al., 2014] Filgueira, R., Atkinson, M., Bell, A., Main, I., Boon, S., Kilburn, C., and Meredith, P. (2014). eScience gateway stimulating collaboration in rock physics and volcanology. In *e-Science (e-Science)*, 2014 IEEE 10th International Conference on, volume 1, pages 187-195. IEEE.
- [Friedland and Iwasaki 1985]: Friedland, P.E.; Iwasaki, Y. J. The Concept and Implementation of Skeletal Plans. *Automatic Reasoning* (1985) 1: 161.
- [Garijo and Gil 2011]: Garijo, D. and Gil, Y. A new approach for publishing workflows: Abstractions, standards, and Linked Data. In *Proceedings of the 6th workshop on Workflows in support of large-scale science*, pages 47-56, Seattle. ACM. 2011.
- [Garijo and Gil 2012]: Garijo, D. and Gil, Y. Augmenting PROV with plans in P-Plan: Scientific processes as Linked Data. In *Second International Workshop on Linked Science: Tackling Big Data (LISC)*, held in conjunction with the *International Semantic Web Conference (ISWC)*, Boston, MA. 2012.
- [Garijo et al., 2013]: Detecting common scientific workflow fragments using templates and execution provenance. In *Proceedings of the Seventh International Conference on Knowledge Capture, K-CAP '13*, pages 33-40. 2013.

- [Garijo et al., 2014]: Garijo, D., Gil, Y., and Corcho, O. Towards workflow ecosystems through semantic and standard representations. In Proceedings of the 9th Workshop on Workflows in Support of Large-Scale Science, WORKS '14, pages 94-104, Piscataway, NJ, 2014.
- [Gil et al., 2007]: Gil, Y.; Deelman, E.; Ellisman, M. H.; Fahringer, T.; Fox, G.; Gannon, D.; Goble, C. A.; Livny, M.; Moreau, L.; and Myers, J. "Examining the Challenges of Scientific Workflows." *IEEE Computer*, 40(12), 2007.
- [Gil et al., 2011]: Gil, Y., Ratnakar, V., Kim, J., Gonzalez-Calero, P. A., Groth, P. T., Moody, J., and Deelman, E. (2011). WINGS: Intelligent workflow-based design of computational experiments. *IEEE Intelligent Systems*, 26(1):62-72.
- [Gil et al., 2011b]: Gil, Y.; Gonzalez-Calero, P. A.; Kim, J.; Moody, J.; and Ratnakar, V. A Semantic Framework for Automatic Generation of Computational Workflows Using Distributed Data and Component Catalogs *Journal of Experimental and Theoretical Artificial Intelligence*, 23(4). 2011.
- [Gil 2013]: Gil, Y. Mapping Semantic Workflows to Alternative Workflow Execution Engines. Gil, Y. In *Seventh IEEE International Conference on Semantic Computing (ICSC)*, Irvine, CA, 2013.
- [Gil et al., 2015]: Gil, Y.; Garijo, D.; Ratnakar, V.; OntoSoft: Capturing Scientific Software Metadata. In *Proceedings of the 8th International Conference on Knowledge Capture*. Palisades, NY. 2015.
- [Gil et al., 2016]: Gil, Y.; Garijo, D.; Ratnakar, V.; Mayani, R.; Adusumilli, R.; Boyce, H.; Mallick, P. Automated Hypothesis Testing with Large Scientific Data Repositories. *Proceedings of the Fourth Annual Conference on Advances in Cognitive Systems*, Evanston, IL, June 2016.
- [Goecks et al., 2010]: Goecks, J., Nekrutenko, A., and Taylor, J. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biology*, 11(8). 2010.
- [Hoekstra and Groth, 2014]: Hoekstra, R. and Groth, P. PROV-O-Viz - Understanding the Role of Activities in Provenance. *Proceedings of the International Provenance and Annotation Workshop*, Cologne. 2014.
- [Heath and Bizer, 2011]: Heath, T. and Bizer, C. *Linked Data: Evolving the Web into a Global Data Space*. Synthesis Lectures on the Semantic Web. Morgan & Claypool Publishers. 2011.
- [Lebo et al., 2013]: Lebo, T., McGuinness, D., Belhajjame, K., Cheney, J., Corsar, D., Garijo, D., Soiland-Reyes, S., Zednik, S., and Zhao, J. The PROV ontology, W3C Recommendation. Technical report, WWW Consortium. 30th April 2013.
- [Leisch 2002]: Leisch, F. "Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis", *Proceedings of Computational Statistics*, 2002.
- [Ludäscher et al., 2006]: Ludäscher, B., Altintas, I., Berkley, C., Higgins, D., Jaeger, E., Jones, M., Lee, E. A., Tao, J., and Zhao, Y. Scientific workflow management and the Kepler system. *Concurrency and Computation: Practice and Experience*, 18(10):1039-1065. 2006.
- [Mattmann et al., 2006]: Mattmann, C. A., Crichton, D. J., Medvidovic, N., and Hughes, S. A software architecture-based framework for highly distributed and data intensive scientific applications. In *Proceedings of the 28th international conference on Software engineering, ICSE '06*, pages 721-730, New York, NY, USA. 2006.
- [Mates et al., 2011]: Mates, P., Santos, E., Freire, J., and Silva, C. T. "CrowdLabs: Social analysis and visualization for the sciences". In *23rd International Conference on Scientific and Statistical Database Management (SSDBM)*, pages 555-564. Springer. 2011.
- [Mesirov 2010]: Mesirov, J. P. "Accessible Reproducible Research." *Science*, 327:415, 2010.

- [Missier et al., 2010]: Missier, P., Sahoo, S. S., Zhao, J., Goble, C., and Sheth, A. (2010). Janus: from Workflows to Semantic Provenance and Linked Open Data. Provenance and Annotation of Data and Processes Third International Provenance and Annotation Workshop IPAW 2010 Troy NY USA June, 2010 Revised Selected Papers 6378, 129-141.
- [Missier et al., 2013]: Missier, P., Dey, S., Belhajjame, K., Cuevas-Vicenttín, V., and Ludäscher, B. (2013). D-PROV: Extending the PROV provenance model with workflow structure. In Proceedings of the 5th USENIX Workshop on the Theory and Practice of Provenance, TaPP '13, pages 9:1-9:7, Berkeley, CA, USA. USENIX Association. 2013.
- [Moreau et al., 2008]: Moreau, L., Ludäscher, B., Altintas, I., Barga, R. S., Bowers, S., Callahan, S., Chin, G., Clifford, B., Cohen, S., Cohen-Boulakia, S., Davidson, S., Deelman, E., Digiampietri, L., Foster, I., Freire, J., Frew, J., Futrelle, J., Gibson, T., Gil, Y., Goble, C., Golbeck, J., Groth, P., Holland, D. A., Jiang, S., Kim, J., Koop, D., Krenek, A., McPhillips, T., Mehta, G., Miles, S., Metzger, D., Munroe, S., Myers, J., Plale, B., Podhorszki, N., Ratnakar, V., Santos, E., Scheidegger, C., Schuchardt, K., Seltzer, M., Simmhan, Y. L., Silva, C., Slaughter, P., Stephan, E., Stevens, R., Turi, D., Vo, H., Wilde, M., Zhao, J. and Zhao, Y. (2008), Special Issue: The First Provenance Challenge. *Concurrency Computat.: Pract. Exper.*, 20: 409–418.
- [Moreau et al., 2011]: Moreau, L., Clifford, B., Freire, J., Futrelle, J., Gil, Y., Groth, P., Kwasnikowska, N., Miles, S., Missier, P., Myers, J., Plale, B., Simmhan, Y., Stephan, E., and den Bussche, J. V. The open provenance model core specification (v1.1). *Future Generation Computer Systems*, 27(6). 2011.
- [Radulovic et al., 2015] Radulovic, F., Poveda-Villalón, M., Vila-Suero, D., Rodriguez-Doncel, V., Garcia-Castro, R., and Gomez-Perez, A. (2015). Guidelines for Linked Data Generation and publication: An example in building energy consumption. *Automation in Construction*, 57:178 - 187.
- [Reich et al., 2006]: Reich, M., Liefeld, T., Gould, J., Lerner, J., Tamayo, P., and Mesirov, J. P. Genepattern 2.0. *Nature genetics*, 38(5):500-501. 2006.
- [Ruiz et al., 2014]: Ruiz, J., Garrido, J., Santander-Vela, J., Sánchez-Expósito, S., and Verdes-Montenegro, L. AstroTaverna: Building workflows with Virtual Observatory services. *Astronomy and Computing*, 7-8:3-11. Special Issue on The Virtual Observatory: I. 2014.
- [Santana-Pérez and Pérez-Hernandez, 2015]: Santana-Pérez, I. and Pérez-Hernandez, M. (2015). Towards reproducibility in scientific workflows: An infrastructure-based approach. *Scientific Programming*, 2015:11.
- [Scheidegger et al., 2008]: Scheidegger, C. E., Vo, H. T., Koop, D., Freire, J., and Silva, C. T. Querying and re-using workflows with VisTrails. In Proceedings of the 2008 ACM SIGMOD international conference on Management of data, SIGMOD '08, pages 1251-1254, New York, USA. ACM. 2008.
- [Shaon et al., 2011]: Shaon, A.; Callaghan, S.; Lawrence, B.; Matthews, B.; Woolf, A.; Osborn, T.; Harpham, C. A Linked Data Approach to Publishing Complex Scientific Workflows. In Proceedings of the IEEE 7th International Conference on eScience, 303-310, Stockholm, 2011.
- [Starlinger et al., 2014]: Starlinger, J., Cohen-Boulakia, S., and Leser, U. Layer decomposition: An effective structure-based approach for scientific workflow similarity. In e-Science (e-Science), 2014 IEEE 10th International Conference on, volume 1, pages 169-176. 2014.
- [Stoyanovich et al., 2010]: Stoyanovich, J., Taskar, B., and Davidson, S. Exploring repositories of scientific workflows. In Proceedings of the 1st International Workshop on Workflow Approaches to New Data-centric Science, Wands '10, pages 7:1-7:10, New York, NY, USA. ACM. 2010.

- [Taylor 2006]: Taylor, I. J. "Triana generations". Presented at: Second IEEE International Conference on e-Science and Grid Computing (e-Science'06), Amsterdam, Netherlands, 4-6 December 2006. E-Science 2006: Second IEEE International Conference on e-Science and Grid Computing: 4-6 December 2006, Amsterdam, Netherlands. Los Alamitos, CA: IEEE, p. 143.
- [Villazon-Terrazas et al. 2011] Villazon-Terrazas, B., Vilches-Blazquez, L., Corcho, O., and Gomez-Perez, A. (2011). Methodological guidelines for publishing government Linked Data. In Wood, D., editor, Linking Government Data, pages 27-49. Springer New York.
- [Wolstencroft et al., 2013]: Wolstencroft, K., Haines, R., Fellows, D., Williams, A., Withers, D., Owen, S., Soiland-Reyes, S., Dunlop, I., Nenadic, A., Fisher, P., Bhagat, J., Belhajjame, K., Bacall, F., Hardisty, A., de la Hidalga, A. N., Vargas, M. P. B., Su, S., and Goble, C. (2013). The Taverna workflow suite: designing and executing workflows of web services on the desktop, web or in the cloud. Nucleic Acids Research.