



Recording provenance of workflow runs with RO-Crate

DOI:

[10.48550/arXiv.2312.07852](https://doi.org/10.48550/arXiv.2312.07852)

Document Version

Submitted manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Leo, S., Crusoe, M. R., Rodríguez-Navas, L., Sirvent, R., Kanitz, A., Geest, P. D., Wittner, R., Pireddu, L., Garijo, D., Fernández, J. M., Colonnelli, I., Gallo, M., Ohta, T., Suetake, H., Capella-Gutierrez, S., Wit, R. D., Kinoshita, B. P., & Soiland-Reyes, S. (2024). Recording provenance of workflow runs with RO-Crate. Manuscript submitted for publication. <https://doi.org/10.48550/arXiv.2312.07852>

Published in:

PLoS ONE

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



Recording provenance of workflow runs with RO-Crate

Simone Leo^{1*}, Michael R. Crusoe^{2,3,4}, Laura Rodríguez-Navas⁵, Raül Sirvent⁵, Alexander Kanitz^{6,7}, Paul De Geest⁸, Rudolf Wittner^{9,10,11}, Luca Pireddu¹, Daniel Garijo¹², José M. Fernández⁵, Iacopo Colonnelli¹³, Matej Gallo⁹, Tazro Ohta^{14,15}, Hirotaka Suetake¹⁶, Salvador Capella-Gutierrez⁵, Renske de Wit², Bruno P. Kinoshita⁵, Stian Soiland-Reyes^{17,18}

1 Center for Advanced Studies, Research, and Development in Sardinia (CRS4), Pula (CA), Italy

2 Vrije Universiteit Amsterdam, Amsterdam, The Netherlands

3 DTL Projects, The Netherlands

4 Forschungszentrum Jülich, Germany

5 Barcelona Supercomputing Center, Barcelona, Spain

6 Biozentrum, University of Basel, Basel, Switzerland

7 Swiss Institute of Bioinformatics, Lausanne, Switzerland

8 VIB Data Core, Gent, Belgium

9 Faculty of Informatics, Masaryk University, Brno, Czech Republic

10 Institute of Computer Science, Masaryk University, Brno, Czech Republic

11 BBMRI-ERIC, Graz, Austria

12 Ontology Engineering Group, Universidad Politécnica de Madrid, Madrid, Spain

13 Computer Science Department, Università degli Studi di Torino, Torino, Italy

14 Database Center for Life Science, Joint Support-Center for Data Science Research, Research Organization of Information and Systems, Shizuoka, Japan

15 Institute for Advanced Academic Research, Chiba University, Chiba, Japan

16 Sator, Incorporated, Tokyo, Japan

17 Department of Computer Science, The University of Manchester, Manchester, United Kingdom

18 Informatics Institute, University of Amsterdam, Amsterdam, The Netherlands

* simone.leo@crs4.it (SL)

Abstract

Recording the provenance of scientific computation results is key to the support of traceability, reproducibility and quality assessment of data products. Several data models have been explored to address this need, providing representations of workflow plans and their executions as well as means of packaging the resulting information for archiving and sharing. However, existing approaches tend to lack interoperable adoption across workflow management systems. In this work we present Workflow Run RO-Crate, an extension of RO-Crate (Research Object Crate) and Schema.org to capture the provenance of the execution of computational workflows at different levels of granularity and bundle together all their associated objects (inputs, outputs, code, etc.). The model is supported by a diverse, open community that runs regular meetings, discussing development, maintenance and adoption aspects. Workflow Run RO-Crate is already implemented by several workflow management systems, allowing interoperable comparisons between workflow runs from heterogeneous systems. We describe the model, its alignment to standards such as W3C PROV, and its implementation in six

workflow systems. Finally, we illustrate the application of Workflow Run RO-Crate in two use cases of machine learning in the digital image analysis domain.

1 Introduction

A crucial part of scientific research is recording the provenance of its outputs. The W3C PROV standard defines provenance as “a record that describes the people, institutions, entities, and activities involved in producing, influencing, or delivering a piece of data or a thing” [1]. Provenance is instrumental to activities such as traceability, reproducibility, accountability, and quality assessment [2]. The constantly growing size and complexity of scientific datasets and the analysis that is required to extract useful information from them has made science increasingly dependent on advanced automated processing techniques in order to get from experimental data to final results [3–5]. Consequently, a large part of the provenance information for scientific outputs consists of descriptions of complex computer-aided data processing steps. This data processing is often expressed as workflows – i.e., high-level applications that coordinate multiple tools and manage intermediate outputs in order to produce the final results.

In order to homogenise the collection and interchange of provenance records, the W3C consortium proposed a standard for representing provenance in the Web (PROV [1]), along with the PROV ontology (PROV-O) [6], an OWL [7] representation of PROV. PROV-O has been widely extended for workflows (e.g., D-PROV [8], ProvONE [9], OPMW [10] (Open Provenance Model for Workflows), P-PLAN [11]), where provenance information is collected in two main forms: prospective and retrospective [12]. *Prospective provenance* – the execution plan – is essentially the workflow itself: it includes a machine-readable specification with the processing steps to be performed and the data and software dependencies to carry out each computation. *Retrospective provenance* refers to what actually happened during an execution – i.e. what were the values of the input parameters, which outputs were produced, which tools were executed, how much time did the execution take, whether the execution was successful or not, etc. Retrospective provenance may be represented at different levels of abstraction, depending on the information that is available and/or required: a workflow execution may be interpreted i) as a single end-to-end activity, ii) as a set of individual execution of workflow steps, or iii) by going a step further and indicating how each step is divided into sub-processes when a workflow is deployed in a cluster. Various workflow management systems, such as WINGS [13] (Workflow INstance Generation and Specialization) and VisTrails [15,16], have adopted PROV and its PROV-O representation to lift the burden of provenance collection from tool users and developers [17,18].

D-PROV, PROV-ONE, OPMW, P-PLAN propose representations of workflow plans and their respective executions, taking into account the features of the workflow systems implementing them (e.g., hierarchical representations, sub-processes, etc.). Other data models, such as *wfprov* and *wfdesc* [19], go a step further by considering not only the link between plans and executions, but also how to package the various artefacts as a Research Object (RO) [20] to improve metadata interoperability and document the context of a digital experiment.

However, while these models address some workflow provenance representation issues, they have two main limitations: first, the extensions of PROV are not directly interoperable because of differences in their granularities or different assumptions in their workflow representations; second, their support from Workflow Management Systems (WMS) is typically one system per model. An early approach to unify and integrate workflow provenance traces across WMSs was the Workflow Ecosystems through STandards (WEST) [14], which used WINGS to build workflow templates and

different converters. In all of these workflow provenance models, the emphasis is on the workflow execution structure as a directed graph, with only partial references for the data items. The REPRODUCE-ME ontology [21] extended PROV and P-PLAN to explain the overall scientific process with the experimental context including real life objects (e.g. instruments, specimens) and human activities (e.g. lab protocols, screening), demonstrating provenance of individual Jupyter Notebook cells [22] and highlighting the need for provenance also where there is no workflow management system.

More recently, interoperability has been partially addressed by Common Workflow Language Prov (CWLProv) [23], which represents workflow enactments as research objects serialised according to the Big Data Bag approach [24]. The resulting format is a folder containing several data and metadata files [25], expanding on the Research Object Bundle approach of Taverna [26]. CWLProv also extends PROV with a representation of executed processes (activities), their inputs and outputs (entities) and their executors (agents), together with their Common Workflow Language (CWL) specification [27] – a standard workflow specification adopted by at least a dozen different workflow systems [28]. Although CWLProv includes prospective provenance as a *plan* within PROV (based on the *wfdesc* model), in practice its implementation does not include tool definitions or file formats. Thus, for CWLProv consumers to reconstruct the full prospective provenance for understanding the workflow, they would also need to inspect the separate workflow definition in the native language of the workflow management system. Additionally, the CWLProv RO may include several other metadata files and PROV serialisations conforming to different formats, complicating its generation and consumption.

As for granularity, CWLProv proposes multiple levels of provenance [23, Figure 2], from Level 0 (capturing workflow definition) to Level 3 (domain-specific annotations). In practice, the CWL reference implementation *cwltool* [30] and the corresponding CWLProv specification [25] record provenance details of all task executions together with the intermediate data and any nested workflows (CWLProv level 2). This level of granularity requires substantial support from the workflow management system implementing the CWL specification, resulting appropriate for workflow languages where the execution plan, including its distribution among the various tasks, is well known in advance. However, it can be at odds with other systems where the execution is more dynamic, depending on the verification of specific runtime conditions, such as the size and distribution of the data (e.g., COMPSs [31]). This design makes the implementation of CWLProv challenging, which the authors suspect may be one of the main causes for the low adoption of CWLProv (at the time of writing the format is supported only by *cwltool*). Finally, being based on the PROV model, CWLProv is highly focused on the interaction between agents, processes and related entities, while support for contextual metadata (such as workflow authors, licence or creation date) in the Research Object Bundle is limited [32] and stored in a separate manifest file, which includes the data identifier mapping to filenames. A project that uses serialised Research Objects similar to those used by CWLProv is Whole Tale [33], a web platform with a focus on the narrative around scientific studies and their reproducibility, where the serialised ROs are used to export data and metadata from the platform. In contrast, our work is primarily focused on the ability to capture the provenance of computational workflow execution including its data and executable workflow definitions.

RO-Crate [34] is an approach for packaging research data together with their metadata and associated resources. RO-Crate extends Schema.org [35], a popular vocabulary for describing resources on the Web. In its simplest form, an RO-Crate is a directory structure that contains a single JSON-LD [36] metadata file at the top level. The metadata file describes all entities stored in the RO-Crate along with their

relationships, and it is both machine-readable and human-readable. RO-Crate is general enough to be able to describe any dataset, but can also be made as specific as needed through the use of extensions called *profiles*. Profiles describe “a set of conventions, types and properties that one minimally can require and expect to be present in that subset of RO-Crates” [100]. The broad set of types and properties from Schema.org, complemented by a few additional terms from other vocabularies, make the RO-Crate model a candidate for expressing a wide range of contextual information that complements and enriches the core information specified by the profile. This information may include, among others, the workflow authors and their affiliations, associated publications, licensing information, related software, etc. This approach is used by WorkflowHub [37], a workflow-system-agnostic workflow registry which specifies a Workflow RO-Crate profile [38] to gather the workflow definition with such metadata in an archived RO-Crate.

In this work, we present **Workflow Run RO-Crate** (WRROC), an extension of RO-Crate for representing computational workflow execution provenance. Our main contributions include:

- a collection of RO-Crate profiles to represent and package both the prospective and the retrospective provenance of a computational workflow run in a way that is machine-actionable [39], independently of the specific workflow language or execution system, and including support for re-execution;
- implementations of this new model in six workflow management systems and in one conversion tool;
- a mapping of our profiles against the W3C PROV-O Standard using the Simple Knowledge Organisation System (SKOS) [40].

To foster usability, the profiles are characterised by different levels of detail, and the set of mandatory metadata items is kept to a minimum in order to ease the implementation. This flexible approach increases the model’s adaptability to the diverse landscape of WMSs used in practice. The base profile, in particular, is applicable to any kind of computational process, not necessarily described in a formal workflow language. All profiles are supported and sustained by the Workflow Run RO-Crate community, which meets regularly to discuss extensions, issues and new implementations.

The rest of this work is organised as follows: we first describe the Workflow Run RO-Crate profiles in Section 2; we then illustrate implementations in Section 3 and usage examples in Section 4; finally, we include a discussion in Section 5 and we conclude the paper with our plans for future work in Section 6.

2 The Workflow Run RO-Crate profiles

RO-Crate profiles are extensions of the base RO-Crate specification that describe how to represent the classes and relationships that appear in a specific domain or use case. An RO-Crate conforming to a profile is not just machine-readable, but also machine-actionable, as a digital object whose type is represented by the profile itself [41].

The Workflow Run RO-Crate profiles are the main outcome of the activities of the Workflow Run RO-Crate Community [42], an open working group that includes workflow users and developers, WMS users and developers, and researchers and software engineers interested in workflow execution provenance and Findable, Accessible, Interoperable and Reusable (FAIR) approaches for data and software. One of the first steps in the development of the Workflow-Run RO-Crate profiles was to compile a list of

requirements to be addressed by the model from all interested participants, in the form of *competency questions* (CQs) [43]. The process also included reviewing existing state of the art models, such as wfprov [19], ProvONE [9] or OPMW [10]. The result was the definition of 11 CQs capturing requirements which span a broad application scope and consider different levels of provenance granularity. Each requirement was supported by a rationale and linked to a GitHub issue to drive the public discussion forward. When a requirement was addressed, related changes were integrated into the profiles and the relevant issue was closed. All the original issues are now closed, and the profiles have had five official releases on Zenodo [50, 51, 55]. The target of several of the original CQs evolved during profile development, as the continuous discussion within the community highlighted the main points to be addressed. This continuous process is reflected in the corresponding issues and pull requests in the community’s GitHub repository. The final implementation of the CQs in the profiles is validated with SPARQL queries that can be run on RO-Crate metadata samples, also available on the GitHub repository [44].

As requirements were being defined, it became apparent that one single profile would not have been sufficient to cater for all possible usage scenarios. In particular, while some use cases required a detailed description of all computations orchestrated by the workflow, others were only concerned with a “black box” representation of the workflow and its execution as a whole (i.e., whether the workflow execution as a whole was successful and which results were obtained). Additionally, some computations involve a data flow across multiple applications that are executed without the aid of a WMS and thus are not formally described in a standard workflow language. These observations led to the development of three profiles:

1. *Process Run Crate*, to describe the execution of one or more tools that contribute to a computation;
2. *Workflow Run Crate*, to describe a computation orchestrated by a predefined workflow;
3. *Provenance Run Crate*, to describe a workflow computation including the internal details of individual step executions.

In the rest of this section we describe each of these profiles in detail. We use the term “class” to refer to a type as defined in RDF(s) and “entity” to refer to an instance of a class. We use italics to denote the properties and classes in each profile: these are defined in the RO-Crate JSON-LD context [46], which extends Schema.org with terms from the Bioschemas [47] ComputationalWorkflow profile [48] and other vocabularies. Note that terms coming from Bioschemas are not specific to the life sciences. We also developed a dedicated term set [49] to represent concepts that are not captured by terms in the RO-Crate context. New terms are defined in RDF(s) following Schema.org guidelines (i.e., using *domainIncludes* and *rangeIncludes* to define domains and ranges of properties). In the rest of the text and images, the following prefixes are used to represent the corresponding namespaces:

s: → <https://schema.org/>
bioschemas: → <https://bioschemas.org/>
bsp: → <https://bioschemas.org/properties/>
wfrun: → <https://w3id.org/ro/terms/workflow-run#>

2.1 Process Run Crate

The Process Run Crate profile [50] contains specifications to describe the execution of one or more software applications that contribute to the same overall computation, but are not necessarily coordinated by a top-level workflow or script (e.g. when executed manually by a human, one after the other as intermediate datasets become available).

The Process Run Crate is the basis for all profiles in the WRROC collection. It specifies how to describe the fundamental classes involved in a computational run: i) a software application represented by a *s:SoftwareApplication*, *s:SoftwareSourceCode* or *bioschemas:ComputationalWorkflow* class; and ii) its execution, represented by a *s:CreateAction* class, and linking to the application via the *s:instrument* property. Other important properties of the *s:CreateAction* class are *s:object*, which links to the action's inputs, and *s:result*, which links to its outputs. The time the execution started and ended can be provided, respectively, via the *s:startTime* and *s:endTime* properties. The *s:Person* or *s:Organization* class that performed the action is specified via the *s:agent* property. Fig 1 shows the classes used in Process Run Crate together with their relationships.

194
195
196
197
198
199
200
201
202
203
204

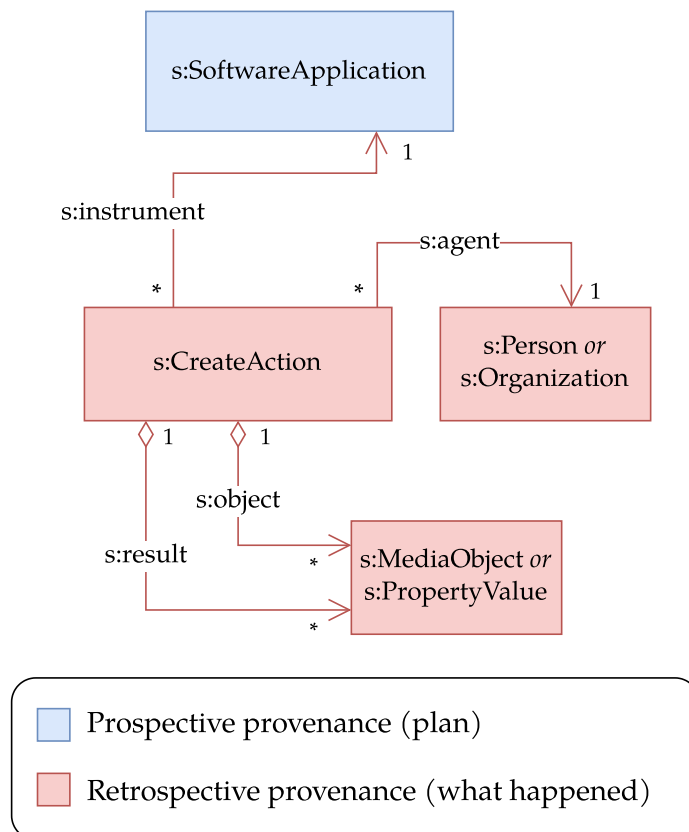


Fig 1. UML class diagram for Process Run Crate. The central class is the *s:CreateAction*, which represents the execution of an application. It links to the application itself via *s:instrument*, to the entity that executed it via *s:agent*, and to its inputs and outputs via *s:object* and *s:result*, respectively. In this and following figures, classes and properties are shown with prefixes to indicate their origin. Some inputs (and, less commonly, outputs) are not stored as files or directories, but passed to the application (e.g., via a command line interface) as values of various types (e.g., a number or string). In this case, the profile recommends a representation via *s:PropertyValue*. For simplicity, we left out the rest of the RO-Crate structure (e.g. the root *s:Dataset*), and attributes (e.g. *s:startTime*, *s:endTime*, *s:description*, *s:actionStatus*). In this UML class notation, diamond \diamond arrows indicate aggregation and regular arrows indicate references, * indicates zero or more occurrences, 1 means single occurrence.

As an example, suppose a user named John Doe runs the UNIX command `head to`

205

extract the first ten lines of an input file named `lines.txt`, storing the result in another file called `selection.txt`. John then runs the `sort` UNIX command on `selection.txt`, storing the sorted output in a new file named `sorted_selection.txt`.

Fig 2 contains a diagram of the two actions and their relationships to the other involved entities. Note how the actions are connected by the fact that the output of “Run Head” is also the input of “Run Sort”: they form an “implicit workflow”, whose steps have been executed manually rather than by a software tool.

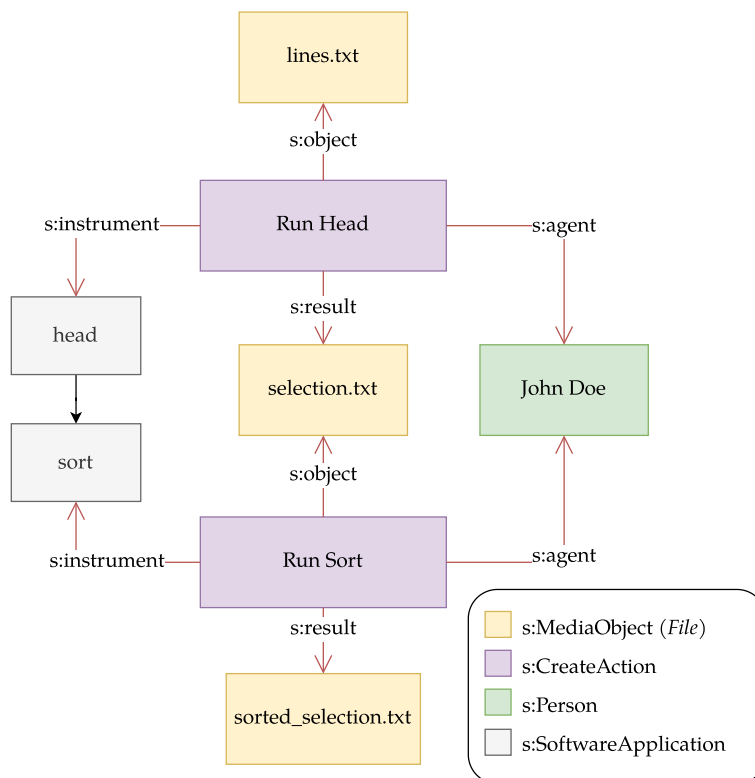


Fig 2. Diagram of a simple workflow where the `head` and `sort` programs were run manually by a user. The executions of the individual software programs are connected by the fact that the file output by `head` was used as input for `sort`, documenting the computational flow in an implicit way. Such executions can be represented with Process Run Crate.

Process Run Crate extends the RO-Crate guidelines on representing software used to create files with additional requirements and conventions. This arrangement is typical of the RO-Crate approach, where the base specification provides general recommendations to allow for high flexibility, while profiles – being more concerned with the representation of specific domains and machine actionability – provide more detailed and structured definitions. Nevertheless, in order to be broadly applicable, profiles also need to avoid the specification of too many strict requirements, trying to strike a good trade-off between flexibility and actionability.

2.2 Workflow Run Crate

The Workflow Run Crate profile [51] combines the Process Run Crate and WorkflowHub’s Workflow RO-Crate [38] profiles to describe the execution of computational workflows managed by a WMS. Such workflows are typically written in a

domain-specific language, such as CWL or Snakemake [52], and run by one or more WMS (e.g., StreamFlow [53], Galaxy [54]). Fig 3 illustrates the classes used in this profile together with their relationships. As in Process Run Crate, the execution is described by a *s:CreateAction* that links to the application via *s:instrument*, but in this case the application must be a workflow, as prescribed by Workflow RO-Crate. More specifically, Workflow RO-Crate states that the RO-Crate must contain a main workflow typed as *File* (an RO-Crate mapping to *s:MediaObject*), *s:SoftwareSourceCode* and *bioschemas:ComputationalWorkflow*. The execution of the individual workflow steps, instead, is not represented: that is left to the more detailed Provenance Run Crate profile (described in the next section).

The Workflow Run Crate profile also contains recommendations on how to represent the workflow’s input and output parameters, based on the Bioschemas ComputationalWorkflow profile. All these elements are represented via the *bioschemas:FormalParameter* class and are referenced from the main workflow via the *bsp:input* and *bsp:output* properties. While the classes referenced from *s:object* and *s:result* in the *s:CreateAction* represent data entities and argument values that were actually used in the workflow execution, the ones referenced from *bsp:input* and *bsp:output* correspond to formal parameters, which acquire a value when the workflow is run (see Fig 3). In the profile, the relationship between an actual value and the corresponding formal parameter is expressed through the *s:exampleOfWork* property. For instance, in the following JSON-LD snippet a formal parameter (*#annotations*) is illustrated together with a corresponding *final-annotations.tsv* file:

```
{
  "@id": "#annotations",
  "@type": "FormalParameter",
  "additionalType": "File",
  "encodingFormat": "text/tab-separated-values",
  "valueRequired": "True",
  "name": "annotations"
},
{
  "@id": "final-annotations.tsv",
  "@type": "File",
  "contentSize": "14784",
  "exampleOfWork": {"@id": "#annotations"}
}
```

2.3 Provenance Run Crate

The Provenance Run Crate profile [55] extends Workflow Run Crate by adding new concepts to describe the internal details of a workflow run, including individual tool executions, intermediate outputs and related parameters. Individual tool executions are represented by additional *s:CreateAction* instances that refer to the tool itself via *s:instrument* – analogously to its use in Process Run Crate. The workflow is required to refer to the tools it orchestrates through the *s:hasPart* property, as suggested in the Bioschemas ComputationalWorkflow profile, though in the latter it is only a recommendation.

To represent the logical steps defined by the workflow, this profile uses *s:HowToStep* – i.e., “A step in the instructions for how to achieve a result” [56]. Steps point to the corresponding tools via the *s:workExample* property and are referenced from the workflow via the *s:step* property; the execution of a step is represented by a *s:ControlAction* pointing to the *s:HowToStep* via *s:instrument* and to the

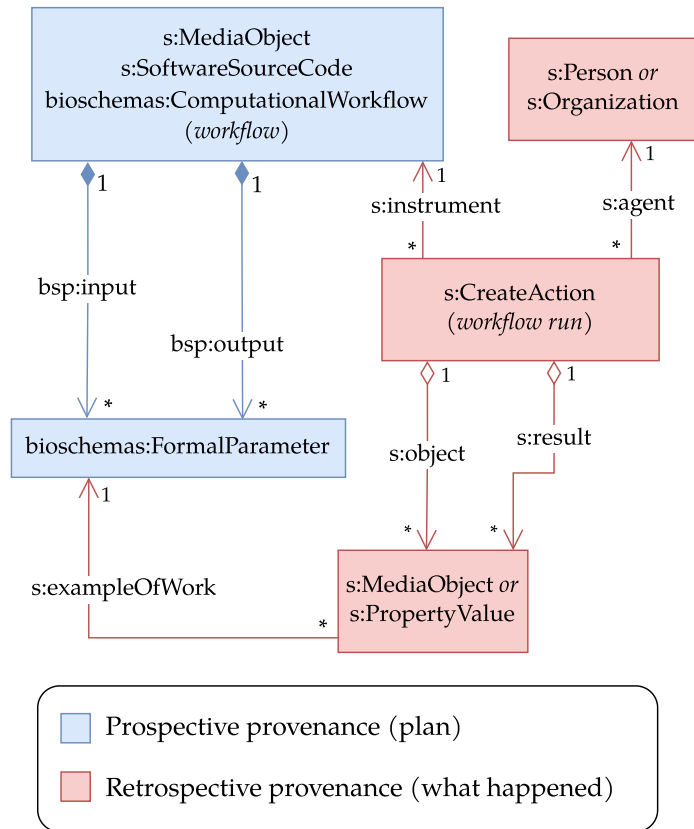


Fig 3. UML class diagram for Workflow Run Crate. The main differences with Process Run Crate are the representation of formal parameters and the fact that the workflow is expected to be an entity with types *s:MediaObject* (*File* in RO-Crate JSON-LD), *s:SoftwareSourceCode* and *bioschemas:ComputationalWorkflow*. Effectively, the workflow belongs to all three types, and its properties are the union of the properties of the individual types. In this profile, the execution history (retrospective provenance) is augmented by a (prospective) workflow definition, giving a high-level overview of the workflow and its input and output parameter definitions (*bioschemas:FormalParameter*). The inner structure of the workflow is not represented in this profile. In the provenance part, individual files (*s:MediaObject*) or arguments (*s:PropertyValue*) are then connected to the parameters they realise. Most workflow systems can consume and produce multiple files, and this mechanism helps to declare each file’s role in the workflow execution. The filled diamond \blacklozenge indicates composition, empty diamond \diamond aggregation, and other arrows relations.

s:CreateAction entities that represent the corresponding tool execution(s) via *s:object*. Note that a step execution does not coincide with a tool execution: an example where this distinction is apparent is when a step maps to multiple executions of the same tool over a list of inputs (e.g. the “scattering” feature in CWL).

An RO-Crate following this profile can also represent the execution of the WMS itself (e.g., cwltool) via *s:OrganizeAction*, pointing to a representation of the WMS via *s:instrument*, to the steps via *s:object* and to the workflow run via *s:result*. The *s:object* attribute of the *s:OrganizeAction* can additionally point to a configuration file containing a description of the settings that affected the behaviour of the WMS during the execution. Fig 4 illustrates the various classes involved in the representation of a workflow run via Provenance Run Crate together with their relationships.

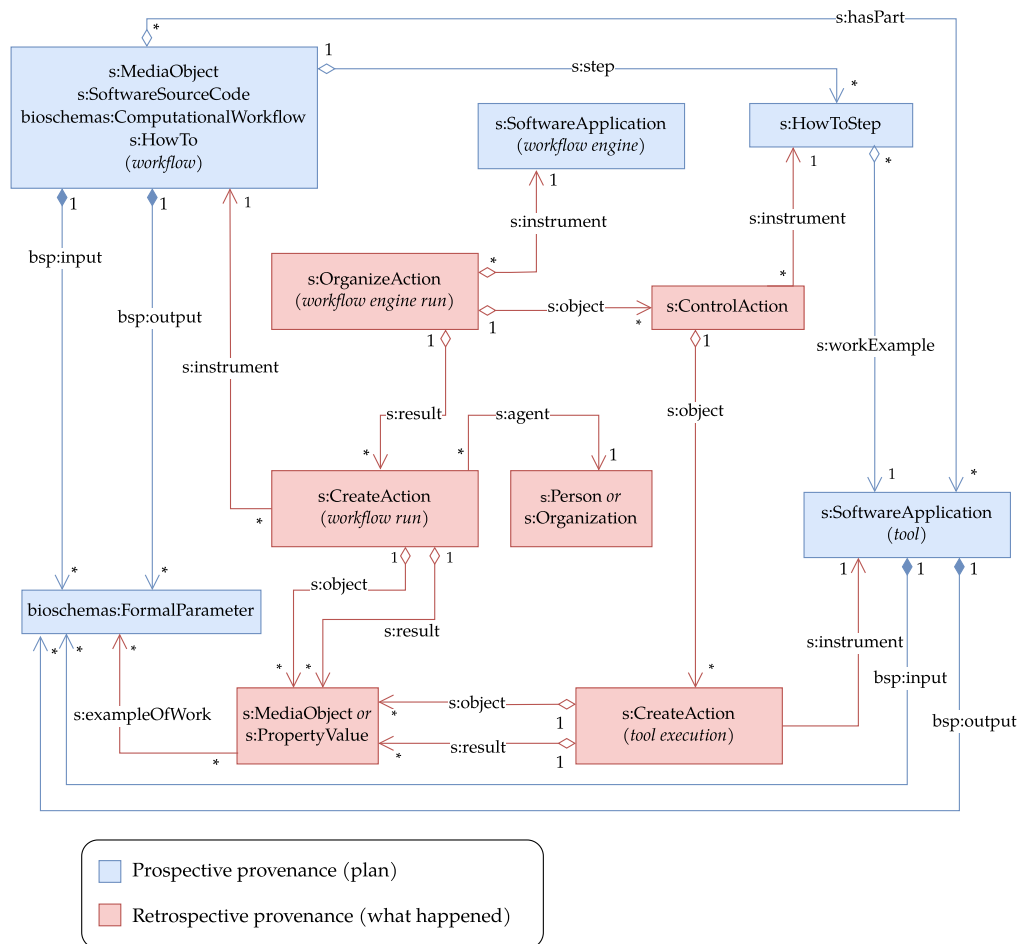


Fig 4. UML class diagram for Provenance Run Crate. In addition to the workflow run, this profile represents the execution of individual steps and their related tools. The prospective side (the execution plan) is shown by the workflow listing a series of *s:HowToSteps*, each linking to the *s:SoftwareApplication* that is to be executed. The *bsp:input* and *bsp:output* parameters for each tool are described in a similar way to the overall workflow parameter in Fig 3. The retrospective provenance side of this profile includes each tool execution as an additional *s:CreateAction* with similar mapping to the realised parameters as *s:MediaObject* or *s:PropertyValue*, allowing intermediate values to be included in the RO-Crate even if they are not workflow outputs. The workflow execution is described the same as in the Workflow Run Crate profile with an overall *s:CreateAction* (the workflow outputs will typically also appear as outputs from inner tool executions). An additional *s:OrganizeAction* represents the workflow engine execution, which orchestrated the steps from the workflow plan through corresponding *s:ControlActions* that spawned the tool’s execution (*s:CreateAction*). It is possible that a single workflow step had multiple such executions (e.g. array iterations). Not shown in figure: *s:actionStatus* and *s:error* to indicate step/workflow execution status. The filled diamond \blacklozenge indicates composition, empty diamond \diamond aggregation, and other arrows relations.

Additionally, this profile specifies how to describe connections between parameters, through *parameter connections* – a fundamental feature of computational workflows. Specifically, parameter connections describe: (i) how tools consume as input the

286
287
288

intermediate outputs generated by other tools; and (ii) how workflow-level parameters are mapped to tool-level parameters. As an example, consider again the workflow depicted in Fig 2, and suppose it is implemented in a workflow language such as CWL: the workflow-level input (a text file) is linked through a parameter connection to the input of the `head` tool wrapper, and then a second parameter connection links this tool's output to the input of the `sort` tool wrapper. A representation of parameter connections is particularly useful for traceability, since it provides the means to document the inputs and tools on which workflow outputs depend. Since the current RO-Crate context has no suitable terms for the description of such relationships, we added appropriate ones to the aforementioned dedicated term set [49]: a *wfrun:ParameterConnection* type with *wfrun:sourceParameter* and *wfrun:targetParameter* attributes that respectively map to the source and target formal parameters, and a *wfrun:connection* property to link from the relevant step or workflow to the *wfrun:ParameterConnection* instances.

In our set of profiles, Provenance Run Crate is the most detailed one and offers the highest level of granularity; its specification is a superset of Workflow Run RO-Crate, which in turn is a superset of Process Run Crate. This relationship between the three profiles is illustrated in Fig 5, as a Venn diagram. Theoretically, all computational provenance information could be represented through the Provenance Run Crate profile alone (possibly relaxing some requirements), since it inherits from the other ones. In practice, though, this choice would require the use of the most complex model even for simple use cases. Having three separate profiles provides a way to represent information at different levels of granularity, while keeping all RO-Crates generated with them interoperable. This approach gives a straightforward path to supporting the representation of computational provenance in simpler use cases such as with simple command executions, i.e. the Process Run Crate. Additionally, the approach lowers the accessibility barrier for implementation in WMSs, as developers may choose to initially implement only the more basic support in their WMS, with reduced effort and complexity, and gradually scale to more detailed representations. This encourages the adoption of WRROC across the diverse landscape of use cases and WMSs.

2.4 Profile formats

The WRROC profiles are available both in human-readable (HTML) and in machine-readable format (RO-Crate). The human-readable profiles are at:

- <https://w3id.org/ro/wfrun/process/0.5>
- <https://w3id.org/ro/wfrun/workflow/0.5>
- <https://w3id.org/ro/wfrun/provenance/0.5>

And the corresponding machine-readable ones at:

- <https://doi.org/10.5281/zenodo.12158562>
- <https://doi.org/10.5281/zenodo.12159311>
- <https://doi.org/10.5281/zenodo.12160782>

The RO-Crate metadata files for the machine readable profiles can be retrieved using the same URLs as the human-readable ones, but with JSON-LD content negotiation: this is done by setting "Accept:application/ld+json" in the HTTP header.

The new terms we defined to represent concepts that could not be expressed with existing Schema.org ones are at:

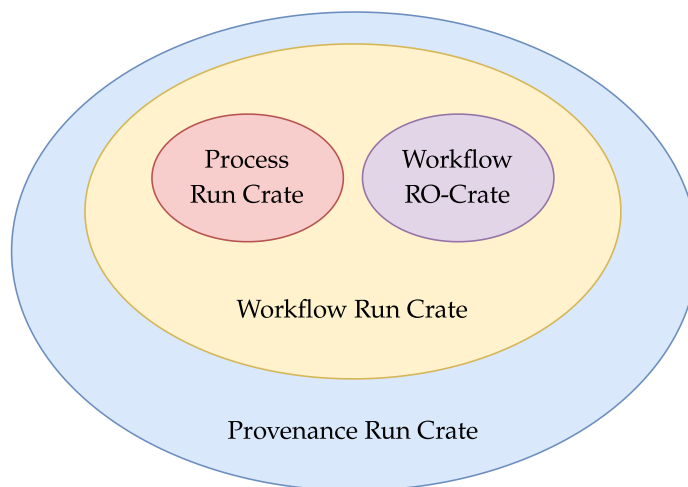


Fig 5. Venn diagram of the specifications for the various RO-Crate profiles. Process Run Crate specifies how to describe the fundamental classes involved in a computational run, and thus is the basis for all profiles in the WRROC collection. Workflow Run Crate inherits the specifications of both Process Run Crate and Workflow RO-Crate. Provenance Run Crate, in turn, inherits the specifications of Workflow Run Crate (and in a sense includes multiple Process Runs for each step execution, but within a single Crate).

- <https://w3id.org/ro/terms/workflow-run>

These terms are available in multiple formats with content negotiation, as explained at the above link.

3 Implementations

Support for the Workflow Run RO-Crate profiles presented in this work has been implemented in a number of systems, showing support from the community and demonstrating their usability in practice. We describe seven of these implementations (one in a conversion tool and six in WMS) in the following sections. Table 1 provides an overview of the implementations, along with the respective profile implemented, and links to the implementation itself and to an example RO-Crate. These tools have been developed in parallel by different teams, and independently from each other. RO-Crate has a strong ecosystem of tools [34], and the WRROC implementations have either re-used these or added their own approach to the standards.

3.1 Runcrate

Runcrate [57] is a Workflow Run RO-Crate toolkit which also serves as a reference implementation of the proposed profiles. It consists of a Python package with a command line interface, providing a straightforward path to integration in Python software and other workflows. The runcrate toolkit includes functionality to convert CWLProv ROs to RO-Crates conforming to the Provenance Run Crate profile (`runcrate convert`), effectively providing an indirect implementation of the format for `cwltool`. Indeed, the CWLProv model provided a basis for the Provenance Run Crate profile, and the implementation of a conversion tool in runcrate at times drove the improvement and extension of the profile as new requirements or gaps in the old designs

emerged. Runcrate converts both the retrospective provenance part of the CWLProv RO (the RDF graph of the workflow’s execution) and the prospective provenance part (the CWL files, including the workflow itself). Both parts are thus converted into a single, workflow-language-agnostic metadata resource.

Another functionality offered by the runcrate package is `runcrate report`, which reports on the various executions described in an input RO-Crate, listing their starting and ending times, the values of the various parameters, etc. Runcrate report demonstrates how the provenance profiles presented in this work enable comparison of runs interoperably across different workflow languages or different implementations of the same language. This functionality has also been used as a lightweight validator for the various implementations.

Runcrate also includes a `run` subcommand to re-execute the computation described by an input Workflow Run Crate or Provenance Run Crate where CWL is used as a workflow language. It works by mapping the RO-Crate description of input parameters and their values (the workflow’s `bsp:input` and the action’s `s:object`) to the format expected by CWL, which is then used to relaunch the workflow on the input data. This functionality shows the machine-actionability of the profiles to support reproducibility, and was used to successfully re-execute the digital pathology annotation workflow described in Section 4.1. Of course, achieving a full re-execution in the general case may not always be possible: reproducibility is supported by the profiles, but also benefits from specific characteristics of the workflow language (which should provide a clear formalism to map input items to their corresponding parameter slots) and of the specific workflow’s implementation, which can be made considerably easier to reproduce by containerising the computational environment required by each step (if allowed by the workflow language).

3.2 Galaxy

The Galaxy project [54] provides a WMS with data management functionalities as a multi-user platform, aiming to make computational biology more accessible to research scientists that do not have computer programming or systems administration experience. Galaxy’s most prominent features include: a collection of 7500+ integrated tools [58]; a web interface that allows the definition and execution of workflows using the integrated tools; a network of dedicated (public) Galaxy instances.

The export of workflow execution provenance data as Workflow Run Crates was added to Galaxy in version 23 [96] providing a more interoperable alternative to the basic export of Galaxy workflow *invocations*. A WRROC export from Galaxy includes: the workflow definition; a set of serialisations of the invocation-related metadata in Galaxy native, JSON-formatted files; and the input and output data files. This result is achieved by: i) extracting provenance data from Galaxy entities related to the workflow run, along with their associated metadata; ii) converting them to RO-Crate metadata using the `ro-crate-py` library [59]; iii) describing all files contained in the basic invocation export within the RO-Crate metadata; and iv) making the Workflow Run Crate available for export to the user through Galaxy’s web interface and API [60]. We extract the prospective provenance contained in Galaxy’s YAML-based `gxformat2` [61] workflow definition, which includes details of the analysis pipeline such as the graph of the tools that need to be executed and metadata about the data types required. The retrospective provenance – i.e., the details of the executed workflow, such as the inputs, outputs, and parameter values used – is extracted from Galaxy’s data model, which is not directly accessible to users in the context of a public Galaxy server. All of this provenance information is then mapped to RO-Crate metadata, including some Galaxy-specific data entities such as dataset collections. An exemplary Workflow Run Crate exported from Galaxy, through its *Workflow Invocations* list, is available on

Zenodo [62].

In practice, a user would take the following steps to obtain a Workflow Run Crate from a Galaxy instance: i) create or download a Galaxy workflow definition (e.g.: from WorkflowHub) and import it in a Galaxy instance, or create a workflow through the Galaxy GUI directly; ii) execute the workflow, providing the required inputs; iii) after the workflow has run successfully, the corresponding RO-Crate will be available for export from the Workflow Invocations list.

3.3 COMPSs

COMPSs [31] is a task-based programming model that allows users to transform a sequential application into a parallel one by simply annotating some of its methods, thus facilitating scaling applications to increasing amounts of computing resources. When a COMPSs application is executed, a corresponding workflow describing the application's tasks and their data dependencies is dynamically generated and used by the COMPSs runtime to orchestrate the execution of the application in the infrastructure. As a WMS, COMPSs stands out for its many advanced features that enable applications to achieve fine-grained high efficiency in HPC systems, such as the ability to exploit underlying parallelisation frameworks (e.g. MPI [63], OpenMP [64]), compilers (e.g. NUMBA [65]), failure management, task grouping, and more. Also, provenance recording for COMPSs workflows has been explored in previous work [66], where the Workflow RO-Crate profile was used to capture structured descriptive metadata about the executed workflow, without introducing any significant run time performance overheads.

In this work, COMPSs has been further improved by implementing the generation of provenance information conformant to the Workflow Run Crate profile, thus also capturing details about the actual execution of the workflow. The dynamic nature of COMPSs workflows poses some challenges to capturing provenance, which were met thanks to the instruments provided by the WRROC model. For instance, a COMPSs workflow is created when the application is executed and, thus, a prior static workflow definition does not exist before that moment. Due to this design, the workflow entity in the metadata file references the entry point application run by COMPSs – instead of, for instance, a dedicated workflow definition file as one might find with other WMSs. Also, formal parameters are not included in the prospective provenance (note that specifying them is not required by the profile) because inputs and outputs (both for each task and the whole workflow) are determined at runtime. However, the RO-Crate generation by COMPSs leverages the information recorded by the runtime to automatically add metadata of all input or output data assets used or produced by the workflow.

Because of the supercomputing environments where COMPSs is used, the integration of Workflow Run Crate support required paying particular attention to the generation of a unique ID for the *s:CreateAction* representing the workflow run. Our implementation uses UUIDs for distributed environments, while it adds a combination of hostname and queuing system job ID for supercomputer executions, to provide as much information as possible from the run while preserving ID uniqueness. In the *s:CreateAction*, the *s:description* term includes system information, as well as relevant environment variables that provide details on the execution environment (e.g., node list, CPUs per node). Finally, the *s:subjectOf* property of the *s:CreateAction* references the system's monitoring tool (when available), where authorised users can see detailed profiling of the corresponding job execution, as provided by the MareNostrum IV supercomputer [67].

To showcase the COMPSs adoption of the Workflow Run Crate profile, we provide as an example the execution of the BackTrackBB [68] application in the MareNostrum IV supercomputer. BackTrackBB targets the detection and location of seismic sources using the statistical coherence of the wave field recorded by seismic networks and antennas. The resulting RO-Crate [69] captures the provenance of the execution results

and complies with the Workflow Run Crate profile. It includes the application source files, a diagram of the workflow’s graph, application profiling and input and output files.

The implementation of provenance recording using Workflow Run Crate has been fully integrated in the COMPSs runtime and is available as of release 3.2 [70].

3.4 StreamFlow

The StreamFlow framework [53] is a container-native WMS for the execution of workflows defined in CWL. It has been designed around two primary principles: first, it allows the execution of tasks in multi-container environments, supporting the concurrent execution of communicating tasks in a multi-agent ecosystem; second, it relaxes the requirement of a single shared data space, allowing for hybrid workflow executions on top of multi-cloud, hybrid cloud/HPC, and federated infrastructures. StreamFlow orchestrates hybrid workflows by combining a *workflow description* (e.g., a CWL workflow description and a set of input values) with one or more *deployment descriptions* – i.e. representations of the execution environments in terms of infrastructure-as-code (e.g., Docker Compose files [71], HPC batch scripts, and Helm charts [72]). A `streamflow.yml` file – the entry point of each StreamFlow execution – binds each workflow step with the set of most suitable execution environments. At execution time, StreamFlow automatically takes care of all the secondary aspects, like scheduling, checkpointing, fault tolerance, and data movements.

StreamFlow collects prospective and retrospective provenance data in a custom format and persists it into a pluggable database (using `sqlite3` as the default choice). After a CWL workflow execution completes, users can generate an RO-Crate through the `streamflow prov` command, which extracts the provenance data stored in the database for one or more workflow executions and converts it to an RO-Crate archive that is fully compliant with the Provenance Run Crate Profile, including the details of each task run by the WMS. Support for the format has been integrated into the main development branch and will be included in release 0.2.0 [73].

From the StreamFlow point of view, the main limitation in the actual version of the Provenance Run Crate standard is the lack of support for distributed provenance – i.e., a standard way to describe the set of locations involved in a workflow execution and their topology. As a temporary solution, the StreamFlow configuration and a description of the hybrid execution environment are preserved by directly including the `streamflow.yml` file into the generated archive. However, this product-specific solution prevents a wider adoption from other WMS and forces agnostic frameworks (e.g., WorkflowHub) to provide ad-hoc plugins to interpret the StreamFlow format. Since the support for hybrid and cross-facility workflows is gaining traction in the WMS ecosystem, we envision support for distributed provenance as a feature for future versions of Workflow Run RO-Crate.

3.5 WfExS-backend

WfExS-backend [74] is a FAIR workflow execution orchestrator that aims to address some of the difficulties found in analysis reproducibility and analysis of sensitive data in a secure manner. WfExS-backend requires that the software used by workflow steps is available in publicly accessible software containers for reproducibility. Actual workflow execution is delegated to one of the supported workflow engines – currently either Nextflow [75] or `cnltool`. The orchestrator prepares and stages all the elements needed to run the workflow – i.e. all the files of the workflow itself, the specific version of the workflow engine, the required software containers and the inputs. All these elements are referenced through resolvable identifiers, ideally public, permanent ones. Thanks to this approach, the orchestrator can consume workflows from various types of sources, such

as git repositories, Software Heritage, or even RO-Crates from WorkflowHub. WfExS-backend development milestones have aimed to reach FAIR workflow execution through the generation and consumption of RO-Crates following the Workflow Run Crate profile, which has proven to be a mechanism suitable to semantically describe digital objects in a way that simplifies embedding details crucial to analysis reproducibility and replicability.

When the orchestrator prepares a workflow for execution it records details relevant to the prospective provenance, such as the public URLs used to fetch input data and workflows, content digestion fingerprints (typically sha256 checksums) and metadata derived from workflow files, container images and input files. Most of this captured metadata is later included in the generated RO-Crates. WfExS-backend has explicit commands to generate and publish both prospective and retrospective provenance RO-Crates based on a given existing staged execution scenario. These RO-Crates can selectively include copies of used elements as payloads. Workflows can be executed more than once in the same staged directory, with all the executions sharing the same inputs. In this case, run details from all the executions are represented in the retrospective provenance RO-Crate. Support for the consumption of Workflow Run RO-Crates to reproduce the operations they document is available as of WfExS-backend version 1.0.0a0 [74]. We have created examples of Workflow Run Crates generated by WfExS-backend to capture provenance information from the execution of a Nextflow workflow [76] and a CWL workflow [30]; these crates are both available on Zenodo [77, 78]. Future developments to WfExS-backend will also add support for embedding in the RO-Crates the URLs of output results that have been deposited into a suitable repository (like Zenodo DOIs, for instance).

3.6 Sapporo

Sapporo [79] is an implementation of the Workflow Execution Service (WES) API specification [109]. WES is a standard proposed by the Global Alliance for Genomics and Health (GA4GH) for cloud-based data analysis platforms that receive requests to execute workflows. Sapporo supports the execution of several workflow engines, including cwltool [30], Toil [80], and StreamFlow [53]. Sapporo includes features specifically tailored to bioinformatics applications, including the calculation of feature statistics from specific types of outputs generated by workflow runs. For example, the system calculates the mapping rate of DNA sequence alignments from BAM format files. To describe the feature values, Sapporo uses the Workflow Run Crate profile extended with additional terms to represent these biological features [81].

Further, the Tonkaz companion command line software has integrated functionality to compare Run Crates generated by Sapporo to measure the reproducibility of the workflow outputs [82]. Developers can use this unique feature to build a CI/CD platform for their workflows to ensure that changes to the product do not produce an unexpected result. Workflow users can also use this feature to verify the results from the same workflow deployed in different environments.

While Sapporo supports Workflow Run Crate, since WES is a WMS wrapper, it does not parse the provided workflow definition files. Instead, it embeds the information in the files passed by the WES request to record the provenance of execution rather than using the actual workflow parameters meant for the wrapped WMS. Therefore, the current implementation of Sapporo does not capture the connections between the inputs/outputs depicted in the workflow and the actual files used/generated during the run. The profile generated by Sapporo has fields representing input and output files, but they are not linked to formal parameters.

Sapporo supports export to Workflow Run Crate as of release 1.5.1 [83]. An example of a Workflow Run RO-Crate generated by Sapporo is available on Zenodo [84].

3.7 Autosubmit

Autosubmit [85] is an open source, lightweight workflow manager and meta-scheduler tailored to configuring and running scientific experiments in climate research. It supports scheduling jobs via SSH to Slurm [86], PBS [87] and other remote batch servers used in HPC.

Autosubmit’s “archive” feature archives the experiment directory and all its contents into a ZIP file, which can be used later to access the provenance data or to execute the Autosubmit experiment again. Even though the data in the ZIP file includes prospective provenance and retrospective provenance, it is not structured, and a simple examination yields no way to distinguish the provenance types.

Recent releases of Autosubmit 4 have added features to increase user flexibility. An updated YAML configuration management system has been implemented that allows users to combine multiple YAML files into a single unified configuration file. Also, the option to use only the experiment manager features of Autosubmit has been added, delegating the workflow execution to a different backend workflow engine – like ecFlow [88], Cylc [89], or a CWL runner. While these features provide some much appreciated flexibility, they have increased the complexity involved in reliably tracking the experiment configuration and other metadata for provenance documentation purposes.

In order to give users a more structured way to archive provenance, which includes the complete experiment configuration, the parameters used to generate it, and is also interoperable between workflow managers, the archive feature was enhanced with a new option in Autosubmit 4.0.100 [90] to enable the generation of provenance data in Workflow Run RO-Crates. The prospective provenance data for the crate is extracted from the Autosubmit experiment configuration. This data includes the multiple YAML files, the unified YAML configuration, as well as the parameters used to preprocess each file – preprocessing replaces placeholders in script templates with values from the experiment configuration. The retrospective provenance data is included with the RO-Crate archive and includes logs and other traces produced by the experiment workflow. Both prospective and retrospective provenance data are included in the final RO-Crate, which is compliant with the Workflow Run Crate profile. At a practical level, the implementation was able to leverage the `ro-crate-py` library for many of the details pertaining to the creation of the RO-Crate archive in Python, and adding information for the JSON-LD metadata.

One of the main challenges for implementing WRROC support in Autosubmit was incorporating Autosubmit’s *Project* feature. A Project in Autosubmit is an abstract concept that references a code repository and is used to define experiment configuration and contains template scripts defining workflow tasks and other auxiliary files. The project has a *type* that defines the *type* of the repository (e.g., git) and a *location* that is the URL to retrieve it. The RO-Crate file generated by Autosubmit includes the project type and location, but it does not include the complete Project and so it is lacking configuration details and scripts. Therefore, users receive provenance data of the Project, but only those with the appropriate privileges can access its constituent resources (many applications run with Autosubmit can not be publicly shared without consent). After consulting with the RO-Crate community regarding the specific Autosubmit requirements, the Autosubmit team adopted a mixed approach where Autosubmit initialises the JSON-LD metadata from its configuration and local trace files, and the user is responsible for providing a partial JSON-LD metadata object in the Autosubmit YAML configuration. `ro-crate-py` was extended to allow the RO-Crate JSON-LD metadata to be patched by these partial JSON-LD metadata objects. This way, users are able to provide the information that is missing from the Autosubmit configuration model, but is required by WRROC – e.g., licence, authors,

inputs, outputs, formal parameters, etc.

Future implementations of WRROC support should be facilitated by the new functionality added to ro-crate-py to support the user-mediated metadata integration approach. On the other hand, the integration of WRROC support would have been facilitated by an automated validation tool for RO-Crate archives, and by documentation and examples on how to use the profiles with *coarse-grained* workflow management systems (as defined in [91]) that do not track inputs and outputs, which is the case of Autosubmit – as well as the Cylc and ecFlow workflow engines. The feedback generated by this use case was welcomed by the WRROC community and work to address these issues is either planned or under way at the time of writing.

To demonstrate Autosubmit’s new WRROC-based functionality to generate structured provenance data, a workflow was created using an example Autosubmit Project designed using UFZ’s mHM (mesoscale Hydrological Model) [93, 94], and it was executed with Autosubmit. The resulting Workflow Run Crate is available from Zenodo [92].

Table 1. Workflow Run Crate implementations

Impl.	Profile	Version URL/DOI	Example
runcrate	Provenance	[57]	[95]
Galaxy	Workflow	[96]	[62]
COMPSs	Workflow	[70]	[69]
Streamflow	Provenance	[73]	[97]
WfExS	Workflow	[74]	[77]
Sapporo	Workflow	[83]	[84]
Autosubmit	Workflow	[90]	[92]

Summary of each WRROC implementation, together with the profile it implements, the software version that makes it available and an example RO-Crate. Runcrate is a toolkit that converts CWLProv ROs to Provenance Run Crates, while the others are workflow management systems.

4 Exemplary use cases

We illustrate Workflow Run RO-Crate on two exemplary use cases. These are similar in terms of application domain, as they both relate to the application of machine learning techniques for the analysis of human prostate images for the purpose of supporting cancer tissue detection. However, the use cases are quite different in the way computations are executed and provenance is represented: in the first, the analysis is conducted by means of a CWL workflow and the outcome is represented with Provenance Run Crate; in the second, Process Run Crate is used in combination with a complementary model to represent a provenance chain that can extend beyond the computational analysis.

4.1 Provenance Run Crate for digital pathology

In this section, we present a use case that demonstrates the effectiveness of the Provenance Run Crate profile at capturing provenance data in the context of digital pathology. More specifically, we demonstrate the generation of RO-Crates to save provenance data associated with the computational annotation of magnified prostate tissue areas and cancer subregions using deep learning models [98]. The image annotation process is implemented in a CWL workflow consisting of three steps, each executing inference on an image using a deep learning model: i) inference of a

low-resolution tissue mask to select areas for further processing; ii) high-resolution tissue inference to refine borders; iii) high-resolution cancer tissue identification. The two tissue inference steps run the same tool, but set different values for the parameter that controls the magnification level, and the second runs on a subset of the image area. The workflow is integrated in the CRS4 Digital Pathology Platform [99], a web-based platform to support clinical studies involving the examination and/or the annotation of digital pathology images.

To assess the interoperability of WRROC, we recorded the provenance of the execution of the same exemplary workflow on two different WMSs. In the first case, we executed the CWL workflow with `cwltool` and converted the resulting CWLProv RO to a Provenance Run Crate with the `runcrate` tool (Section 3.1). In the second case, the workflow was executed with the StreamFlow WMS (Section 3.4). The RO-Crates obtained in the two cases [95,97] are very similar to each other, differing only in a few details. For instance, Streamflow includes its configuration file in the crate and has separate files for the workflow and the two tools, while `cwltool` with `runcrate` results in the workflow and the tools being stored in a single file (CWL’s “packed” format). Apart from these minor differences, the description of the computation is essentially the same, so the RO-Crates are fully interoperable. Four actions are represented: the workflow itself, the two executions of the tissue extraction tool and the execution of the tumour classification tool. Each action is linked to the corresponding workflow or tool via the *s:instrument* property, and reports its starting and ending time. For each action, input and output slots are referenced by the workflow, while the corresponding values are referenced by the action itself. The data and *s:PropertyValue* entities corresponding to the input and output values link to the corresponding parameter slots via the *s:exampleOfWork* property, providing information on the values taken by the parameters during execution. Listing 1 shows the output of the `runcrate report` command for the StreamFlow RO-Crate. For each action (workflow or tool run), `runcrate` reports the associated instrument (workflow or tool), the starting and ending time and the list of inputs and outputs, with pointers from the formal parameter to the corresponding actual value taken during the execution of the action.

The *s:exampleOfWork* link between input / output values and parameter slots is used by `runcrate run` to reconstruct the CWL input parameter mapping needed to rerun the computation. The *s:alternateName* property (a Schema.org property applicable to all entities), which records the original name of data entities (at the time the computation was run), is also crucial for reproducibility in this case: both StreamFlow and CWLProv, to avoid clashes, record input and output files and directories using their SHA1 checksum as their names. However, for this particular workflow file names are important: it expects the input image data to be in the MIRAX [101] format, where the “main” dataset file taken as an input parameter by the processing application must be accompanied by a directory of additional data files, in the same location and with the same name, apart from the extension. The `runcrate` tool uses the *s:alternateName* to rename the input dataset as required, so that the expected pattern can be picked up by the workflow during the re-execution. This use case was the main motivation to include a recommendation to use *s:alternateName* with the above semantics in Process Run Crate.

Thanks to the fact that both RO-Crates were generated following the best practices to support reproducibility mentioned in the profiles, we were able to automatically re-execute both computations with the `runcrate` tool. This was also made possible by the fact that the CWL workflow included information on which container images to use for each tool. Overall, this shows how reproducibility is a hard-to-achieve goal that can only be supported, but not ensured, by the profiles, since it also depends on factors like the characteristics of the computation, the choice of workflow language and whether

Listing 1. Output of the `runcrate report` command executed on the Provenance Run Crate generated by StreamFlow in the digital pathology inference use case (Section 4.1). This informal listing of relevant RO-Crate entities describes each step of the execution. Note that inputs and outputs are of different types (not shown): e.g., `tissue_low>0.9` is a string parameter, `6b15de...` is a filename, and `#af0253...` is a collection.

```

action: #30a65cba-1b75-47dc-ad47-1d33819cf156
  instrument: predictions.cwl (['SoftwareSourceCode',
    'ComputationalWorkflow', 'HowTo', 'File'])
  started: 2023-05-09T05:10:53.937305+00:00
  ended: 2023-05-09T05:11:07.521396+00:00
  inputs:
    #af0253d688f3409a2c6d24bf6b35df7c4e271292 <- predictions.cwl#slide
    tissue_low <- predictions.cwl#tissue-low-label
    9 <- predictions.cwl#tissue-low-level
    tissue_low>0.9 <- predictions.cwl#tissue-high-filter
    tissue_high <- predictions.cwl#tissue-high-label
    4 <- predictions.cwl#tissue-high-level
    tissue_low>0.99 <- predictions.cwl#tumor-filter
    tumor <- predictions.cwl#tumor-label
    1 <- predictions.cwl#tumor-level
  outputs:
    06133ec5f8973ec3cc5281e5df56421c3228c221 <- predictions.cwl#tissue
    4fd6110ee3c544182027f82ffe84b5ae7db5fb81 <- predictions.cwl#tumor
action: #457c80d0-75e8-46d6-bada-b3fe82ea0ef1
  step: predictions.cwl#extract-tissue-low
  instrument: extract_tissue.cwl (['SoftwareApplication', 'File'])
  started: 2023-05-09T05:10:55.236742+00:00
  ended: 2023-05-09T05:10:55.910025+00:00
  inputs:
    tissue_low <- extract_tissue.cwl#label
    9 <- extract_tissue.cwl#level
    #af0253d688f3409a2c6d24bf6b35df7c4e271292 <- extract_tissue.cwl#src
  outputs:
    6b15de40dd0ee3234062d0f261c77575a60de0f2 <- extract_tissue.cwl#tissue
action: #d09a8355-1a14-4ea4-b00b-122e010e5cc9
  step: predictions.cwl#extract-tissue-high
  instrument: extract_tissue.cwl (['SoftwareApplication', 'File'])
  started: 2023-05-09T05:10:58.417760+00:00
  ended: 2023-05-09T05:11:03.153912+00:00
  inputs:
    tissue_low>0.9 <- extract_tissue.cwl#filter
    6b15de40dd0ee3234062d0f261c77575a60de0f2 <- extract_tissue.cwl#filter_slide
    tissue_high <- extract_tissue.cwl#label
    4 <- extract_tissue.cwl#level
    #af0253d688f3409a2c6d24bf6b35df7c4e271292 <- extract_tissue.cwl#src
  outputs:
    06133ec5f8973ec3cc5281e5df56421c3228c221 <- extract_tissue.cwl#tissue
action: #ae2163a8-1a2a-4d78-9c81-caad76a72e47
  step: predictions.cwl#classify-tumor
  instrument: classify_tumor.cwl (['SoftwareApplication', 'File'])
  started: 2023-05-09T05:10:58.420654+00:00
  ended: 2023-05-09T05:11:06.708344+00:00
  inputs:
    tissue_low>0.99 <- classify_tumor.cwl#filter
    6b15de40dd0ee3234062d0f261c77575a60de0f2 <- classify_tumor.cwl#filter_slide
    tumor <- classify_tumor.cwl#label
    1 <- classify_tumor.cwl#level
    #af0253d688f3409a2c6d24bf6b35df7c4e271292 <- classify_tumor.cwl#src
  outputs:
    4fd6110ee3c544182027f82ffe84b5ae7db5fb81 <- classify_tumor.cwl#tumor

```

best practices such as containerisation are followed.

This use case highlighted the need to add specifications on how to represent multi-file datasets [50, section “Representing multi-file objects”], driven by the need to handle the aforementioned MIRAX image format. To represent these, we added specifications to the Process Run Crate profile on describing “composite” datasets consisting of multiple files and directories to be treated as a single unit – as opposed to more conventional input or output parameters consisting of a single file. The profile specifies that such datasets should be represented by a *s:Collection* class linking to individual files and directories via the *s:hasPart* property, and referencing the main part (if any) via the *s:mainEntity* property. Note that, by adding this specification to Process Run Crate, we also made it available to Workflow Run Crate and Provenance Run Crate. In the output of the `runcrate` report tool the additional files are not shown, since the formal parameter points to the *s:Collection* class that describes the whole dataset.

This use case also demonstrates the usage of parameter connections (described in Section 2.3). The RO-Crate resulting from the workflow run contains a representation of all connections between workflow-level parameters (the overall input and output parameters) and tool-level parameters. This allows crate consumers to programmatically find which tool is affected by a workflow-level parameter, thus providing insight on how the workflow works internally (the main feature of the Provenance Run Crate profile). For instance, the `tissue-high-level` workflow parameter is connected to the `level` parameter of the `extract_tissue.cwl` tool by the `extract-tissue-high` step. This parameter regulates the resolution level (pyramidal images are organised into multiple levels of resolution) at which the image is processed in the high-resolution tissue extraction phase. A similar connection is present for the tissue extraction at low resolution. Since *ufrun:ParameterConnections* are referenced from the relevant *s:HowToStep*, the crate consumer can easily determine the resolution level used for both image processing phases from the retrospective provenance.

4.2 Process Run Crate and CPM RO-Crate for cancer detection

This section presents an RO-Crate created to describe an execution of a computational pipeline that trains AI models to detect the presence of carcinoma cells in high-resolution digital images of magnified human prostate tissue. This RO-Crate makes use of Process Run Crate and CPM RO-Crate [102], an RO-Crate profile that supports the representation of entities described according to the Common Provenance Model (CPM) [103, 104, 106].

The CPM is a recently developed extension of the W3C PROV model [1]. It enables the representation of distributed provenance, which is created when an object involved in the research process – either digital or physical (e.g., biological material) – is exchanged between organisations, so that each organisation can document only a portion of the object’s life cycle. Using CPM, each involved organisation can document its portion of the life cycle by generating, storing, and managing individual provenance components, which are then linked together in a chain that spans multiple organizations. The CPM prescribes how to represent such provenance, and how to enable its traversal and processing using a common algorithm, independently from the type of object being described. In addition, the CPM defines a notion of meta-provenance, which contains metadata about the history of individual provenance components.

CPM RO-Crate supports the identification of CPM-based provenance and meta-provenance files within an RO-Crate, so that data, metadata, and CPM-based provenance information can be packed together. An RO-Crate generated according to the CPM-RO-Crate profile embeds parts of the distributed provenance, which may be linked to the provenance of precursors and successors of the packed data. The CPM-RO-Crate profile synergises well with Process Run Crate, since the former can

add references to CPM-based provenance descriptions of computational executions described with the latter, integrating them in the distributed provenance. Since CPM-based provenance and meta-provenance files are typically themselves produced by computations, Process Run Crate allows to represent these along with the main computations that produce the datasets being exchanged, providing the full picture in a cohesive ensemble.

The use case pipeline consists of three main computational steps: i) a preprocessing step that splits input images into small patches and divides them into a training and a testing set; ii) a training step that trains the model to recognise the presence of carcinoma cells in the images; iii) an evaluation step that measures the accuracy of the trained model on the testing set. In addition to these pipeline steps, the RO-Crate describes additional computations related to the generation of the CPM provenance and meta-provenance files. All computations are described according to the Process Run Crate profile, while the CPM files are referenced according to the CPM RO-Crate profile. Also represented via Process Run Crate are: the input dataset; the results of the pipeline execution; the scripts that implement the pipeline; the log files generated by the scripts; a script that converts the logs into the CPM files. This approach allowed us to describe all elements as a single RO-Crate, which is available on Zenodo [105].

Listing 2 presents the `runcrate report` output for the RO-Crate, including action inputs and outputs while omitting other details. The listing shows the connections between the actions, forming an “implicit workflow” as discussed in Section 2.1. For instance, the `prov_train.log` file is both an output of the training action (`#train_script:ROCRATE-PUB-...`) and an input of the CPM provenance generation action for the training phase (`#train_script:6efa9a06-...:CPM-provgen`), highlighting the interdependency between the steps.

The CPM files complement the RO-Crate with details about the pipeline execution process, such as how the input dataset was split into training and testing sets, or detailed information about each training iteration of the AI model. For instance, the RO-Crate contains a representation of a checkpoint of the AI model after the second training iteration, with the corresponding entity’s attributes containing paths to the respective model stored as a file. The entity is related to the respective training iteration activity, which contains the iteration parameters represented as an attribute list. In addition, the CPM generally provides means to link the input dataset provenance to the provenance of its precursors – human prostate tissues and biological samples the tissues were derived from; this is not included in the example because we used a publicly available input database for which provenance of the precursors was not available. However, the linking mechanism for provenance precursors is exactly the same as between the bundles for the AI pipeline parts. While the RO-Crate is focused on the execution of the pipeline, the provenance included in the CPM files intends to be interlinked with provenance of the precursors or successors, providing means to traverse the whole provenance chain. For the described digital pathology pipeline, the precursors would be: i) a biological sample acquired from a patient; ii) slices of the sample processed and put on glass slides; iii) the images created as a result of scanning the slides using a microscope. As a result, combining the CPM and RO-Crate enables the lookup of research artefacts related to the computation across heterogeneous organisations using the underlying provenance chain.

5 Discussion

The RO-Crate profiles presented in this work provide a unified data model to describe the prospective and retrospective provenance of the execution of a computational workflow, together with contextual metadata about the workflow itself and its

Listing 2. Excerpt of the output of the `runcrate report` command for the AI model training Process Run Crate; only inputs and outputs of the actions are shown. The listing shows the connections between the pipeline actions through the entities they produce or consume – e.g., `cam16_mrxs.h5` is output of the conversion script `convert_script:ff67...` and input for the training script `train_script:ROCRATE...`

```
action: #convert_script:ff67ce65-736f-46d5-9fec-10953cad8695
  inputs:
    wsi/test/
    wsi/train/
    prov_converter_config.json
  outputs:
    cam16_mrxs.h5
    prov_preprocess.log

action: #test_script:ROCRATE-PUB-1438b57a750ce887d4433d9e
  inputs:
    prov_test_config.json
    cam16_mrxs.h5
  outputs:
    predictions.h5
    prov_test.log

action: #test_script:d3cfd9cf-6851-43c6-bee9-c8dc18f22368:CPM-provgen
  inputs:
    prov_test.log
  outputs:
    prov_test.provn
    prov_test.provn.log
    prov_test.png

action: #train_script:ROCRATE-PUB-1438b57a750ce887d4433d9e
  inputs:
    prov_train_config.json
    cam16_mrxs.h5
  outputs:
    prov_train.log
    model/weights/auc_01.ckpt.index
    model/weights/auc_01.ckpt.data-00000-of-00001
    model/weights/auc_02.ckpt.index
    model/weights/auc_02.ckpt.data-00000-of-00001
    model/weights/best_loss.ckpt.index
    model/weights/best_loss.ckpt.data-00000-of-00001
    model/weights/auc_03.ckpt.index
    model/weights/auc_03.ckpt.data-00000-of-00001

action: #train_script:6efa9a06-b8e9-4cfc-88c7-e9d35e5263c3:CPM-provgen
  inputs:
    prov_train.log
  outputs:
    prov_train.provn
    prov_train.png
    prov_train.provn.log

action: #convert_script:9d030b68-70d8-4526-82fe-160d9cfe4806:CPM-provgen
  inputs:
    prov_preprocess.log
  outputs:
    prov_preprocess.provn
    prov_preprocess.png
    prov_preprocess.provn.log

action: #meta_provn_script:86bae258-4c51-4215-854b-32cb49f239ab:CPM-provgen
  inputs:
    prov_train.provn.log
    prov_test.provn.log
    prov_preprocess.provn.log
  outputs:
    meta_provenance.provn
    meta_provenance.png
    meta_provenance.provn.log
```


associated entities (inputs, outputs, code, etc.). The profiles are flexible, allowing one to tailor the provenance description to a broad variety of use cases, agnostic to the WMS used, and allow describing provenance traces at different levels of granularity. These characteristics facilitate implementing support in workflow systems. Six WMS have already integrated support for a WRROC profile, as described in Section 3. These new RO-Crate profiles enable interoperability between implementations, which has been demonstrated through the comparison of workflow executions on heterogeneous systems.

Choosing to base our approach on the RO-Crate model has led to a number of benefits. The collected provenance data can be treated with standard RDF tools. As an example, the following SPARQL [107] query returns all actions in a Workflow Run RO-Crate, together with their instruments and their starting and ending times, independently of the original workflow type or the WMS that executed the workflow:

```
PREFIX schema: <https://schema.org/>
SELECT ?action ?instrument ?start ?end
WHERE {
  ?action a schema:CreateAction .
  ?action schema:instrument ?instrument .
  OPTIONAL { ?action schema:startTime ?start } .
  OPTIONAL { ?action schema:endTime ?end }
}
```

Further, having workflow runs and plans described according to the RO-Crate model allows capturing the context of the workflow itself (e.g. authors, related publications, other workflows, etc.), in addition to the trace alone. Another advantage of RO-Crate is that the files corresponding to the data entities (inputs, outputs, code, etc.) do not necessarily have to be stored together with the metadata file: for instance, they can be remote and referred to via an http(s) URI. This aspect is mostly relevant in situations where the file is very large or cannot be shared publicly, since a URI can reference a resource to which access is limited (e.g., accessible only after authentication, or from specific network boundaries, etc.).

The WRROC profiles are extensions of the base RO-Crate specification that specialise it for the use case of workflow execution provenance representation. The additional terms, constraints and recommendations introduced by the profiles allow users to represent classes and relationships involved in a workflow execution in a precise and detailed way, so that consumers of the RO-Crate can programmatically retrieve the relevant information according to predefined patterns and act upon it. This is a crucial advantage over using the base RO-Crate specification, which was not designed to answer the competency questions defined for capturing the provenance of workflow executions.

The ability to build FAIR into Workflow Management Systems was identified as one of the current open challenges in the Scientific Workflows domain at the Workflows Community Summit [108], with the objective of achieving FAIR Computational Workflows. The profiles introduced in this article help tackle this challenge by introducing interoperable metadata among WMSs that captures the provenance of their corresponding workflow executions. The derivation of Workflow Run Crate, and in turn Provenance Run Crate, from Workflow RO-Crate makes the digital objects that conform to these new profiles compatible with the WorkflowHub workflow registry [37]. This design entails that Workflow Run RO-Crates directly reference the workflow with which the provenance was generated, and it allows workflow runs to be registered on WorkflowHub and easily found and shared with other researchers. Additionally, the inheritance mechanism allows reusing the specifications already developed for Workflow RO-Crate, which form part of the guidelines on representing the prospective provenance.

The Workflow Run RO-Crate profiles, the associated tooling, the implementations and the examples are developed and supported by the open WRROC Community. At

the time of writing, the Community numbers nearly 40 members and brings together members of the RO-Crate community [34], WMS users and developers, workflow users and developers, GA4GH [109] Cloud developers and provenance model authors, and is open to anyone who is interested in the representation of workflow execution provenance. The inclusion of WMS developers and workflow users has been key to keeping the specifications flexible, easy to implement and grounded on real use cases, while the diversity of the stakeholders has included a plurality of viewpoints while driving the model’s development forward, resulting in profiles that are already being used (as described in Section 3).

In the following subsections, we provide an evaluation of the metadata coverage of runcrate and we discuss how WRROC relates to standards such as W3C PROV-O and to other community projects.

5.1 Evaluation of metadata coverage using runcrate convert

Since CWLProv was a starting point in the development of WRROC (Section 3.1), as a baseline validation we chose to verify that the metadata contained in CWLProv ROs is preserved in the RO-Crates produced by their conversion through runcrate’s *convert* command. In previous work we had conducted a qualitative analysis of metadata coverage in CWLProv (version 0.6.0), based on concrete examples of ROs associated with a realistic bioinformatics workflow [110]; in this work we repeated this analysis for WRROC, and compared the WRROC RDF representation (in `ro-crate-metadata.json`) with the CWLProv RDF provenance graph. To summarise, the analysis focuses on the comparison of the degree of representation by the two models of six provenance data types defined in [110], which we recall here for clarity.

- T1. Scientific context:** the choices which were made in the design of the workflow and parameter values.
- T2. Data:** input and output data.
- T3. Software:** the tools directly orchestrated by the workflow, and their dependencies.
- T4. Workflow:** the workflow and tool descriptions, but not the software they control.
- T5. Computational environment:** metadata about the system on which the workflow was executed, comprising both software and hardware.
- T6. Execution details:** additional information about the workflow execution itself.

Each type is in turn articulated in a set of data subtypes, forming a hierarchy of elements that should be represented in workflow provenance data to satisfy a range of use cases spanning from supporting workflow development to supporting a service based on the execution of the workflow, with several other use cases in between. For a full motivation and description of the criteria the reader may refer to the original work [110].

Our analysis shows that, overall, most of the information contained in the CWLProv RDF is transferred to the RO-Crate metadata. The results are summarised in Table 2; for completeness, we also report the (non-RDF) representation of provenance metadata in CWL-specific documents (`packed.cwl` and `primary-job.json`), which are included in both CWLProv ROs and RO-Crates generated by runcrate. We observe that out of the total 20 provenance data subtypes that are part of the analysis, WRROC represented 13 (65%) of them (9 fully, 4 partially), while CWLProv RDF captured 8 (3 fully, 5 partially). The representation of some entire categories of metadata has improved – notably Workflow parameters (WF2), which were insufficiently described in

CWLProv RDF, but defined with type and format in RO-Crate. Moreover, the Workflow Run RO-Crate RDF contains a representation of tools orchestrated by the workflow (T3), as well as a much more extensive description of the workflow itself (T4) compared to CWLProv.

In conclusion, our analysis shows that `runcrate` preserves most provenance metadata previously shown to be relevant in realistic RO use case scenarios. More detailed results of the analysis can be found in [111].

Table 2. Summarised results of our qualitative analysis of Provenance Run Crates generated with `runcrate`.

CWL (non-RDF)	Type	Subtype	Name	CWLProv RDF	WRROC RDF
•	T1	SC1	Workflow design	.	•
○		SC2	Entity annotations	.	.
.		SC3	Workflow execution ann.	.	.
○	T2	D1	Data identification	.	.
○		D2	File characteristics	○	○
○		D3	Data access	.	.
•		D4	Parameter mapping	•	•
•	T3	SW1	Software identification	.	•
•		SW2	Software documentation	.	•
•		SW3	Software access	.	•
•	T4	WF1	Workflow software	○	•
•		WF2	Workflow parameters	○	•
•		WF3	Workflow requirements	.	○
.	T5	ENV1	Software environment	.	.
.		ENV2	Hardware environment	.	.
○		ENV3	Container image	○	○
.	T6	EX1	Execution timestamps	•	•
.		EX2	Consumed resources	.	.
.		EX3	Workflow engine	○	○
.		EX4	Human agent	•	•

We converted CWLProv (v0.6.0) ROs to WRROC with `runcrate` 0.5.0. The table compares the degree to which the data subtypes of the provenance data taxonomy (identified by the triple (Type, Subtype, Name)) are preserved by the CWLProv RDF and the WRROC RDF models; the taxonomy is defined in previous work [110], where relevant provenance metadata are identified based on realistic use cases for ROs associated with a real-life bioinformatics workflow. For completeness, the *CWL (non-RDF)* column also reports the non-RDF representation of provenance metadata in CWL-specific documents: `packed.cwl` (the workflow) and `primary-job.json` (the input parameter file). Since `packed.cwl` and `primary-job.json` are also included in RO-Crate, we only considered how the metadata was represented in `ro-crate-metadata.json`.

Legend: • fully represented ○ partially represented . missing or unstructured representation

5.2 Workflow Run RO-Crate and the W3C PROV standard

One of our aims for the WRROC profiles is to make them compatible with both Schema.org and W3C PROV. Provenance Run Crate is the profile that most closely matches the level of detail provided by CWLProv, which extends W3C PROV. Table 3 shows how the main classes and relationships represented by Provenance Run Crate map to PROV constructs, using the SKOS vocabulary to indicate the type of relationship between each pair of terms. A machine-readable version of the mapping can be found in the RO-Crate accompanying this article [112, 113].

Table 3. Mapping from Workflow Run RO-Crate to equivalent W3C PROV concepts using SKOS [40]. For instance, *s:CreateAction* has **broader** match *prov:Activity*, meaning that *prov:Activity* is more general. Prefix *prov:* <https://www.w3.org/ns/prov#>.

RO-Crate	Relationship	W3C PROV-O
<i>s:Action</i> (superclass of <i>s:CreateAction</i> , <i>s:OrganizeAction</i>)	Has close match (Schema.org Actions may also be potential actions in the future)	<i>prov:Activity</i>
<i>s:CreateAction</i> , <i>s:OrganizeAction</i>	Has broader match	<i>prov:Activity</i>
<i>s:Person</i>	Has exact match	<i>prov:Person</i>
<i>s:Organization</i>	Has exact match	<i>prov:Organization</i>
<i>s:SoftwareApplication</i>	Has related match	<i>prov:SoftwareAgent</i>
<i>bioschemas:ComputationalWorkflow</i> , <i>s:SoftwareApplication</i> , <i>s:HowTo</i>	Has broader match	<i>prov:Plan</i> , <i>prov:Entity</i>
<i>s:MediaObject</i> , <i>s:Dataset</i> , <i>s:PropertyValue</i>	Has broader match	<i>prov:Entity</i>
<i>s:startTime</i> on <i>s:CreateAction</i>	Has close match	<i>prov:startedAtTime</i>
<i>s:endTime</i> on <i>s:CreateAction</i>	Has close match	<i>prov:endedAtTime</i>
<i>s:agent</i> on <i>s:CreateAction</i>	Has related match	<i>prov:wasStartedBy</i> , <i>prov:wasEndedBy</i>
<i>s:agent</i> and <i>s:instrument</i> on <i>s:CreateAction</i>	Has broader match	<i>prov:wasAssociatedWith</i>
<i>s:instrument</i> on <i>s:CreateAction</i>	Has related match (Complex mapping: an instrument implies a qualified association with the agent, linked to a plan)	<i>prov:hadPlan</i> on <i>prov:Association</i>
<i>s:object</i> on <i>s:CreateAction</i>	Has exact match	<i>prov:used</i>
<i>s:result</i> on <i>s:CreateAction</i>	Has close match	inverse <i>prov:wasGeneratedBy</i>

5.3 Five Safes Workflow Run Crate

The *Five Safes RO-Crate* [114] profile has been developed to extend the Workflow Run Crate profile for use in Trusted Research Environments (TRE), following the Five Safes Framework [116] to better handle sensitive health data in federated workflow execution across TREs in the UK [115]. A crate with a workflow run request references a pre-approved workflow and project details for manual and automated assessment according to the TRE’s agreement policy for the sensitive dataset. The crate then goes through multiple phases internal to the TRE, including validation, sign-off, workflow execution and disclosure control. At this stage the crate is also conforming to the Workflow Run Crate profile. The final crate is then safe to be made public.

This extension of Workflow Run Crate documents and supports the *human review process* – important for transparency on TRE data usage. The initial implementation of this process used WfExS as the workflow execution backend, and this approach will form the basis for further work on implementing federated workflow execution in the British initiatives DARE UK and HDR UK [117] and in the European EOOSC-ENTRUST project for Trusted Research Environments [118].

5.4 Biocompute Object RO-Crate

IEEE 2791-2020 [119], colloquially known as *Biocompute Objects* (BCO), is a standard for representing provenance of a genomic sequencing pipeline, intended for submission of the workflow to regulatory bodies – e.g. as part of a personalised medical treatment method [120]. The BCO is represented as a single JSON file which includes description of the workflow and its steps and intended purpose, as well as references for tools used and data sources accessed. There is overlap in the goals of BCO and Workflow Run Crate profiles, however their intentions and focus are different. BCO is primarily conveying a computational method for the purpose of manual regulatory review and further reuse, with any values provided as an exemplar run. A Workflow Run Crate, however, is primarily documenting a particular workflow execution, and the workflow is associated to facilitate rerun rather than reuse.

Previously, a guide to packaging BioCompute Objects using RO-Crate [121] was developed as a profile to combine both standards [122]. In this early approach, RO-Crate was primarily a vessel to transport the BCO along with its constituent resources, including the workflow and data files, as well as to provide these resources with additional typing and licence metadata that is not captured by the BCO JSON. Further work is being planned with the BCO community to update the BCO-RO profile to align with the newer Workflow Run RO-Crate profiles.

6 Conclusion and future work

The Workflow Run RO-Crate profile collection presented in this manuscript is a new model to represent and package both the prospective and the retrospective provenance relating to the execution of computational workflows in a way that is machine-actionable, interoperable, independent of the specific workflow language or execution system, and including support for re-execution. These new profiles build on RO-Crate and Schema.org to include contextual information and bundle together all objects of the workflow execution (inputs, outputs, code, etc.). Our approach minimizes the set of mandatory metadata items and defines a hierarchy of profiles – Process Run Crate, Workflow Run Crate, and Provenance Run Crate – that capture provenance information at increasing levels of detail and complexity. This flexible approach increases the model’s adaptability to the diverse landscape of WMSs used in practice, and modulates the implementation effort as a function of the requirements of the specific use case. As a result, there has already been significant uptake of Workflow Run RO-Crate, as shown by its adoption in six WMS, including Galaxy, StreamFlow and COMPSs; in addition, the `runcrate` toolkit has been implemented as part of this work providing various inspection, conversion and re-execution functionalities. Moreover, we have shown how WRROC has been applied in real use cases.

Workflow Run RO-Crate is an ongoing project. Therefore, our profiles and the surrounding software are not static entities, but keep being updated to cater for new requirements and use cases. As examples of ongoing work, at the time of writing there are plans to expand the `runcrate` toolkit to better support the creation and querying of WRROC objects. Also, work is ongoing to implement automated conformance validation of crates. In addition, several of the implementations presented in this work will also develop new features. For instance, the Galaxy community plans to extend its WRROC support to: include metadata detailing each step of a workflow run to conform to the Provenance Run Crate profile; develop and/or integrate RO-Crate more deeply with import and export of Galaxy histories; and further develop user-guided import of RO-Crates as Galaxy datasets, histories and workflows. Further, we are currently exploring the cloud execution of Workflow Run RO-Crates. The Workflow Execution

Service (WES) specification is used by the Global Alliance for Genomics and Health (GA4GH) [109] to enable WMS-agnostic interpretation of workflows and scheduling of task execution. In addition, the Task Execution Service (TES) specification enables the execution of individual, atomic, containerised tasks in a compute backend-independent manner. We are planning to undertake an in-depth analysis of the degree of interoperability between the TES and WES API standards – roughly the equivalents of Process and Workflow Run Crates, respectively – by placing their focus on the actual execution of tasks/processes and workflows in cloud environments and liaising with the GA4GH Cloud community to align schemas where necessary. We will then build an interconversion library that attempts to i) construct WES workflow and TES task run requests from RO-Crates containing Provenance, Workflow or Process Run requests and therefore allow their easy (re)execution on any GA4GH Cloud API-powered infrastructure, and ii) bundle information from the WES and TES (as well as other GA4GH Cloud API resources, where available) to create or extend RO-Crates with standards-compliant Process, Workflow or even Provenance RO-Crates.

The maintenance and development of WRROC is driven by an open community, currently numbering about 40 members. The Community runs regular virtual meetings (every two weeks at the time of writing) and coordinates on Slack and the RO-Crate mailing list. Naturally, feedback and contributions from the community are welcome and encouraged, and new requirements and features are discussed and sustained, particularly through the WRROC GitHub repository issue tracker [45]. Through the open Community we expect to encourage and support further adoption of WRROC, be it by the other WMS or other use cases, maybe in time converging towards a common workflow execution provenance representation.

Acknowledgments

The authors would like to thank all participants to the Workflow Run RO-Crate working group meetings for the fruitful discussions and valuable feedback.

Author contributions

Author contributions following the CRediT Taxonomy:

Simone Leo Conceptualization, Data Curation, Investigation, Methodology, Resources, Software, Supervision, Validation, Visualization, Writing – Original Draft preparation, Writing – Review & Editing

Michael R. Crusoe Conceptualization, Investigation, Software, Supervision

Laura Rodríguez-Navas Software, Writing – Original Draft preparation

Raül Sirvent Data Curation, Software, Writing – Original Draft preparation, Writing – Review & Editing

Alexander Kanitz Writing – Original Draft preparation, Writing – Review & Editing

Paul De Geest Data Curation, Software, Writing – Original Draft preparation

Rudolf Wittner Data Curation, Writing – Original Draft preparation, Writing – Review & Editing

Luca Pireddu Funding acquisition, Project Administration, Supervision, Writing – Review & Editing

Daniel Garijo Conceptualization, Formal Analysis, Writing – Original Draft preparation, Writing – Review & Editing	1017 1018
José M. Fernández Data Curation, Software, Writing – Original Draft preparation	1019
Iacopo Colonnelli Data Curation, Software, Writing – Original Draft preparation	1020
Matej Gallo Data Curation, Software	1021
Tazro Ohta Data Curation, Software, Writing – Original Draft preparation	1022
Hiroataka Suetake Data Curation, Software, Writing – Original Draft preparation	1023
Salvador Capella-Gutierrez Funding Acquisition, Resources, Supervision, Writing – Original Draft preparation	1024 1025
Renske de Wit Software, Writing – Original Draft preparation, Writing – Review & Editing	1026 1027
Bruno de Paula Kinoshita Data Curation, Software, Writing – Original Draft preparation, Writing – Review & Editing	1028 1029
Stian Soiland-Reyes Conceptualization, Formal Analysis, Funding Acquisition, Investigation, Methodology, Resources, Software, Supervision, Visualization, Writing – Original Draft preparation, Writing – Review & Editing	1030 1031 1032

References

1. Moreau L, Missier P, Belhajjame K, B'Far R, Cheney J, Coppens S, et al. PROV-DM: The PROV Data Model. W3C Recommendation 30 April 2013 [cited 2023 Dec 7]. <https://www.w3.org/TR/2013/REC-prov-dm-20130430/>
2. Herschel M, Diestelkämper R, Ben Lahmar H. A survey on provenance: What for? What form? What from? The VLDB Journal, 2017;26:881–906. doi: [10.1007/s00778-017-0486-1](https://doi.org/10.1007/s00778-017-0486-1)
3. Himanen L, Geurts A, Foster AS, Rinke P. Data-Driven Materials Science: Status, Challenges, and Perspectives. Advanced Science, 2019;6(21):1900808. doi: [10.1002/advs.201900808](https://doi.org/10.1002/advs.201900808)
4. Gauthier J, Vincent AT, Charette SJ, Derome N. A brief history of bioinformatics. Briefings in Bioinformatics, 2019;20(6):1981–1996. doi: [10.1093/bib/bby063](https://doi.org/10.1093/bib/bby063)
5. Huntingford C, Jeffers ES, Bonsall MB, Christensen HM, Lees T, Yang H. Machine learning and artificial intelligence to aid climate change research and preparedness. Environmental Research Letters, 2019;14(12):124007. doi: [10.1088/1748-9326/ab4e55](https://doi.org/10.1088/1748-9326/ab4e55)
6. Lebo T, Sahoo S, McGuinness D, Belhajjame K, Cheney J, Corsar D, et al. PROV-O: The PROV Ontology. W3C Recommendation 30 April 2013 [cited 2023 Dec 7]. <https://www.w3.org/TR/2013/REC-prov-o-20130430/>
7. W3C OWL Working Group. OWL 2 Web Ontology Language Document Overview (Second Edition). W3C Recommendation 11 December 2012 [cited 2023 Dec 7]. <http://www.w3.org/TR/2012/REC-owl2-overview-20121211/>
8. Missier P, Dey S, Belhajjame K, Cuevas-Vicenttín V, Ludäscher B. 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), 2013.
9. Cuevas-Vicenttín V, Ludäscher B, Missier P, Belhajjame K, Chirigati F, Wei Y, et al. ProvONE: A PROV Extension Data Model for Scientific Workflow Provenance, 2016 [cited 2023 Dec 7]. <https://purl.dataone.org/provone-v1-dev>

10. Garijo D, 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 (WORKS '11) 2011. doi: [10.1145/2110497.2110504](https://doi.org/10.1145/2110497.2110504)
11. Garijo D, Gil Y. Augmenting PROV with Plans in P-PLAN: Scientific Processes as Linked Data. In Proceedings of the Second International Workshop on Linked Science, 2012.
12. Freire J, Koop D, Santos E, Silva CT. Provenance for Computational Tasks: A Survey. *Computing in Science & Engineering* 2012;10(3):11–21. doi: [10.1109/MCSE.2008.79](https://doi.org/10.1109/MCSE.2008.79)
13. Gil Y, Ratnakar V, Kim J, Gonzalez-Calero P, Groth P, Moody J, et al. Wings: Intelligent Workflow-Based Design of Computational Experiments. *IEEE Intelligent Systems* 2011;26(1). doi: [10.1109/MIS.2010.9](https://doi.org/10.1109/MIS.2010.9)
14. Garijo D, Gil Y, Corcho O. Towards Workflow Ecosystems through Semantic and Standard Representations. In Proceedings of the 9th Workshop on Workflows in Support of Large-Scale Science 2014. doi: [10.1109/works.2014.13](https://doi.org/10.1109/works.2014.13)
15. Scheidegger CE, Vo HT, Koop D, Freire J, Silva CT. Querying and re-using workflows with VisTrails. In Proceedings of the 2008 ACM SIGMOD international conference on Management of data 2008. doi: [10.1145/1376616.1376747](https://doi.org/10.1145/1376616.1376747)
16. Costa F, Silva V, de Oliveira D, Ocaña K, Ogasawara E, Dias J, et al. Capturing and querying workflow runtime provenance with PROV: a practical approach. In Proceedings of the Joint EDBT/ICDT 2013 Workshops 2013. doi: [10.1145/2457317.2457365](https://doi.org/10.1145/2457317.2457365)
17. Atkinson M, Gesing S, Montagnat J, Taylor I. Scientific workflows: Past, present and future. *Future Generation Computer Systems* 2017;75:216–227. doi: [10.1016/j.future.2017.05.041](https://doi.org/10.1016/j.future.2017.05.041)
18. Pérez B, Rubio J, Sáenz-Adán C. A systematic review of provenance systems. *Knowledge and Information Systems* 2018;57:495–543. doi: [10.1007/s10115-018-1164-3](https://doi.org/10.1007/s10115-018-1164-3)
19. Belhajjame K, Zhao J, Garijo D, Gamble M, Hettne K, Palma R, et al. Using a suite of ontologies for preserving workflow-centric research objects. *Journal of Web Semantics* 2015;32:16–42. doi: [10.1016/j.websem.2015.01.003](https://doi.org/10.1016/j.websem.2015.01.003)
20. Bechhofer S, Buchan I, De Roure D, Missier P, Ainsworth J, Bhagat J, et al. Why linked data is not enough for scientists. *Future Generation Computer Systems* 2013;29(2):599–611. doi: [10.1016/j.future.2011.08.004](https://doi.org/10.1016/j.future.2011.08.004)
21. Samuel S, König-Ries B. End-to-End provenance representation for the understandability and reproducibility of scientific experiments using a semantic approach. *Journal of Biomedical Semantics* 2022;13:1. doi: [10.1186/s13326-021-00253-1](https://doi.org/10.1186/s13326-021-00253-1)
22. Samuel S, König-Ries B. ProvBook: Provenance-based Semantic Enrichment of Interactive Notebooks for Reproducibility. The 17th International Semantic Web Conference (ISWC) 2018 Demo Track, 2018.
23. Khan FZ, Soiland-Reyes S, Sinnott RO, Lonie A, Goble C, Crusoe MR. Sharing interoperable workflow provenance: A review of best practices and their practical application in CWLProv. *GigaScience* 2019;8(11):giz095. doi: [10.1093/gigascience/giz095](https://doi.org/10.1093/gigascience/giz095)
24. Chard K, D'Arcy M, Heavner B, Foster I, Kesselman C, Madduri R, et al. I'll take that to go: Big data bags and minimal identifiers for exchange of large, complex datasets. 2016 IEEE International Conference on Big Data (Big Data) 2016;319–328. doi: [10.1109/BigData.2016.7840618](https://doi.org/10.1109/BigData.2016.7840618)
25. Soiland-Reyes S, Khan FZ, Crusoe MR. common-workflow-language/cwlprov: CWLProv 0.6.0. Zenodo, 2018. doi: [10.5281/zenodo.1471585](https://doi.org/10.5281/zenodo.1471585)
26. Soiland-Reyes S, Alper P, Goble C. Tracking workflow execution with TavernaProv. Zenodo, 2016. doi: [10.5281/zenodo.51314](https://doi.org/10.5281/zenodo.51314)

27. Crusoe MR, Abeln S, Iosup A, Amstutz P, Chilton J, Tijanić N, et al. Methods Included: Standardizing Computational Reuse and Portability with the Common Workflow Language. *Communications of the ACM*, 2022;65(6):54–63. doi: [10.1145/3486897](https://doi.org/10.1145/3486897)
28. Common Workflow Language Implementations [cited 2024 May 24]. <https://www.commonwl.org/implementations/>
29. Soiland-Reyes S. The Roterms ontology. Release 30 July 2015 [cited 2024 May 24]. <https://wf4ever.github.io/ro/2016-01-28/roterms/>
30. Amstutz P, Crusoe MR, Khan FZ, Soiland-Reyes S, Singh M, Kumar K, et al. common-workflow-language/cwltool: 3.1.20230127121939. Zenodo, 2023. doi: [10.5281/zenodo.7575947](https://doi.org/10.5281/zenodo.7575947)
31. Lordan F, Tejedor E, Ejarque J, Rafanell R, Álvarez J, Marozzo F, et al. ServiceSs: An interoperable programming framework for the cloud. *Journal of Grid Computing* 2014;12:67–91. doi: [10.1007/s10723-013-9272-5](https://doi.org/10.1007/s10723-013-9272-5)
32. Research Object Bundle context [cited 2024 May 24] <https://w3id.org/bundle/context>
33. Chard K, Gaffney N, Jones MB, Kowalik K, Ludäscher B, McPhillips T, et al. Application of BagIt-Serialized Research Object Bundles for Packaging and Re-Execution of Computational Analyses. 2019 15th International Conference on eScience (eScience) 2019. doi: [10.1109/eScience.2019.00068](https://doi.org/10.1109/eScience.2019.00068)
34. Soiland-Reyes S, Sefton P, Crosas M, Castro LJ, Coppens F, Fernández JM, et al. Packaging research artefacts with RO-Crate. *Data Science* 2022;5(2):97–138. doi: [10.3233/DS-210053](https://doi.org/10.3233/DS-210053)
35. Guha RV, Brickley D, Macbeth S. Schema.org: Evolution of Structured Data on the Web: Big data makes common schemas even more necessary. *Queue* 2015;13(9):10–37. doi: [doi:10.1145/2857274.2857276](https://doi.org/10.1145/2857274.2857276)
36. Sporny M, Longley D, Kellogg G, Lanthaler M, Champin PA, Lindström N. JSON-LD 1.1: A JSON-based Serialization for Linked Data. W3C Recommendation 16 July 2020 [cited 2023 Dec 11]. <https://www.w3.org/TR/2020/REC-json-ld11-20200716/>
37. Goble C, Soiland-Reyes S, Bacall F, Owen S, Williams A, Eguinoa I, et al. Implementing FAIR Digital Objects in the EOSC-Life Workflow Collaboratory. Zenodo, 2021. doi: [10.5281/zenodo.4605654](https://doi.org/10.5281/zenodo.4605654)
38. Bacall F, Williams AR, Owen S, Soiland-Reyes S. Workflow RO-Crate Profile 1.0. WorkflowHub community, 2022 [cited 2023 Dec 11]. <https://w3id.org/workflowhub/workflow-ro-crate/1.0>
39. Batista D, Gonzalez-Beltran A, Sansone SA, Rocca-Serra P. Machine actionable metadata models. *Scientific Data*, 2022;9:592. doi: [10.1038/s41597-022-01707-6](https://doi.org/10.1038/s41597-022-01707-6)
40. Isaac A, Summers E. SKOS Simple Knowledge Organization System Primer. W3C Working Group Note 18 August 2009 [cited 2023 Dec 11]. <https://www.w3.org/TR/2009/NOTE-skos-primer-20090818/>
41. Soiland-Reyes S, Sefton P, Castro LJ, Coppens F, Garijo D, Leo S, et al. Creating lightweight FAIR Digital Objects with RO-Crate. *Research Ideas and Outcomes*, 2022;8:e93937. doi: [10.3897/rio.8.e93937](https://doi.org/10.3897/rio.8.e93937)
42. Workflow Run RO-Crate [cited 2024 May 24]. <https://www.researchobject.org/workflow-run-crate>
43. Workflow Run RO-Crate competency questions [cited 2024 May 24]. <https://www.researchobject.org/workflow-run-crate/requirements>
44. SPARQL queries for the Competency Questions [cited 2024 June 4]. <https://github.com/ResearchObject/workflow-run-crate/tree/main/docs/sparql>

45. Workflow Run RO-Crate GitHub repository [cited 2024 July 2].
<https://github.com/ResearchObject/workflow-run-crate>
46. RO-Crate JSON-LD context, version 1.1 [cited 2024 May 24].
<https://www.researchobject.org/ro-crate/1.1/context.jsonld>
47. Gray A, Goble C, Jimenez R, The Bioschemas Community (2017). Bioschemas: From Potato Salad to Protein Annotation. ISWC (Posters, Demos & Industry Tracks), 2017.
<https://iswc2017.semanticweb.org/paper-579/>
48. Bioschemas ComputationalWorkflow Profile, version 1.0-RELEASE (09 March 2021) [cited 2024 May 24].
<https://bioschemas.org/profiles/ComputationalWorkflow/1.0-RELEASE>
49. ro-terms: Workflow run namespace [cited 2024 Jul 03].
<https://w3id.org/ro/terms/workflow-run>
50. Workflow Run RO-Crate working group. Process Run Crate specification. Version 0.5. Zenodo, 2024. doi: [10.5281/zenodo.12158562](https://doi.org/10.5281/zenodo.12158562)
51. Workflow Run RO-Crate working group. Workflow Run Crate specification. Version 0.5. Zenodo, 2024. doi: [10.5281/zenodo.12159311](https://doi.org/10.5281/zenodo.12159311)
52. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* 2012;28(19):2520–2522. doi: [10.1093/bioinformatics/bts480](https://doi.org/10.1093/bioinformatics/bts480)
53. Colonnelli I, Cantalupo B, Merelli I, Aldinucci M. StreamFlow: cross-breeding Cloud with HPC. *IEEE Transactions on Emerging Topics in Computing*, 2021;9(4):1723–1737. doi: [10.1109/TETC.2020.3019202](https://doi.org/10.1109/TETC.2020.3019202)
54. The Galaxy Community. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2022 update. *Nucleic Acids Research* 2022;50(W1):W345–W351. doi: [10.1093/nar/gkac247](https://doi.org/10.1093/nar/gkac247)
55. Workflow Run RO-Crate working group. Provenance Run Crate specification. Version 0.5. Zenodo, 2024. doi: [10.5281/zenodo.12160782](https://doi.org/10.5281/zenodo.12160782)
56. Schema.org HowToStep definition [cited 2024 May 24]. <https://schema.org/HowToStep>
57. Leo S, Soiland-Reyes S, Crusoe MR. Runcrate. Version 0.5.0. Zenodo, 2023. doi: [10.5281/zenodo.10203433](https://doi.org/10.5281/zenodo.10203433)
58. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, et al. Dissemination of scientific software with Galaxy ToolShed. *Genome Biology* 2014;15:403. doi: [10.1186/gb4161](https://doi.org/10.1186/gb4161)
59. De Geest P, Driesbeke B, Eguinoa I, Gaignard A, Huber S, Kinoshita B, et al. ResearchObject/ro-crate-py: ro-crate-py 0.9.0. Zenodo, 2023. doi: [10.5281/zenodo.10017862](https://doi.org/10.5281/zenodo.10017862)
60. De Geest P, Coppens F, Soiland-Reyes S, Eguinoa I, Leo S. Enhancing RDM in Galaxy by integrating RO-Crate. *Research Ideas and Outcomes*, 2022;8:e95164. doi: [10.3897/rio.8.e95164](https://doi.org/10.3897/rio.8.e95164)
61. Galaxy Workflow Format 2 Description [cited 2024 May 24].
https://galaxyproject.github.io/gxformat2/v19_09.html
62. De Geest P. Run of an example Galaxy collection workflow. Zenodo, 2023. doi: [10.5281/zenodo.7785861](https://doi.org/10.5281/zenodo.7785861)
63. Gabriel E, Fagg GE, Bosilca G, Angskun T, Dongarra JJ, Squyres JM et al. Open MPI: Goals, Concept, and Design of a Next Generation MPI Implementation. *Lecture Notes in Computer Science*, 2004;3241:97–104. doi: [10.1007/978-3-540-30218-6_19](https://doi.org/10.1007/978-3-540-30218-6_19).

64. Dagum L, Menon R. OpenMP: an industry standard API for shared-memory programming. *IEEE Computational Science and Engineering* 1998;5(1):46-55. doi: [10.1109/99.660313](https://doi.org/10.1109/99.660313).
65. Lam SK, Pitrou A, Seibert S. Numba: a LLVM-based Python JIT compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC* 2015. doi: [10.1145/2833157.2833162](https://doi.org/10.1145/2833157.2833162).
66. Sirvent R, Conejero J, Lordan F, Ejarque J, Rodriguez-Navas L, Fernandez JM, et al. Automatic, Efficient, and Scalable Provenance Registration for FAIR HPC Workflows. *2022 IEEE/ACM Workshop on Workflows in Support of Large-Scale Science (WORKS)*, 2022. doi: [10.1109/works56498.2022.00006](https://doi.org/10.1109/works56498.2022.00006)
67. MareNostrum 4 user's guide [cited 2024 May 24]. <https://bsc.es/supportkc/docs/MareNostrum4/intro/>
68. Poiata N, Satriano C, Vilotte JP, Bernard P, Obara K. Multiband array detection and location of seismic sources recorded by dense seismic networks. *Geophysical Journal International*, 2016;205(3):1548–1573. doi: [10.1093/gji/ggw071](https://doi.org/10.1093/gji/ggw071)
69. Poiata N, Satriano C, Conejero J. BackTrackBB: Multi-band array detection and location of seismic sources (PyCOMPSs implementation). Zenodo, 2023. doi: [10.5281/zenodo.7788030](https://doi.org/10.5281/zenodo.7788030)
70. Ejarque J, Lordan F, Badia RM, Sirvent R, Lezzi D, Vazquez F, et al. COMPSs. Version v3.2. Zenodo, 2023. doi: [10.5281/zenodo.7975340](https://doi.org/10.5281/zenodo.7975340)
71. Reis D, Piedade B, Correia FF, Dias JP, Aguiar A. Developing Docker and Docker-Compose Specifications: A Developers' Survey. *IEEE Access*, 2022;10:2318–2329. doi: [10.1109/ACCESS.2021.3137671](https://doi.org/10.1109/ACCESS.2021.3137671)
72. Zerouali A, Opdebeeck R, De Roover C. Helm Charts for Kubernetes Applications: Evolution, Outdatedness and Security Risks. *2023 IEEE/ACM 20th International Conference on Mining Software Repositories*, 2023;523–533. doi: [10.1109/MSR59073.2023.00078](https://doi.org/10.1109/MSR59073.2023.00078)
73. Colonnelli I, Cantalupo B, Aldinucci M, Saitta G, Mulone A. StreamFlow. Version 0.2.0.dev10. Software Heritage Archive, 2023. <https://identifiers.org/swh:1:rev:b2014add57189900fa5a0a0403b7ae3a384df73b>
74. Fernández JM, Rodríguez-Navas L, Muñoz-Cívico A, Iborra P, Lea D. WfExS-backend. Version 1.0.0a0. Zenodo, 2024. doi: [10.5281/zenodo.1258912](https://doi.org/10.5281/zenodo.1258912)
75. Di Tommaso P, Chatzou M, Floden EW, Prieto Barja P, Palumbo E, Notredame C. Nextflow enables reproducible computational workflows. *Nature Biotechnology* 2017;35:316–319. doi: [10.1038/nbt.3820](https://doi.org/10.1038/nbt.3820)
76. Bouyssié D, Altmer P, Capella-Gutierrez S, Fernández JM, Hagemeyer YP, Horvatovich P, et al. WOMBAT-P: Benchmarking Label-Free Proteomics Data Analysis Workflows. *Journal of Proteome Research*, 2023. doi: [10.1021/acs.jproteome.3c00636](https://doi.org/10.1021/acs.jproteome.3c00636)
77. Fernández González JM. RO-Crate from staged WfExS working directory 047b6dfc-3547-4e09-92f8-df7143038ff4 (overbridging templon). Zenodo, 2024. doi: [10.5281/zenodo.12588049](https://doi.org/10.5281/zenodo.12588049)
78. Fernández JM. RO-Crate from staged WfExS working directory a37fee9e-4288-4a9e-b493-993a867207d0 (meer oxometalate). Zenodo, 2024. doi: [10.5281/zenodo.12622362](https://doi.org/10.5281/zenodo.12622362)
79. Suetake H, Tanjo T, Ishii M, Kinoshita BP, Fujino T, Hachiya T, et al. Sapporo: A workflow execution service that encourages the reuse of workflows in various languages in bioinformatics [version 1; peer review: 2 approved with reservations]. *F1000Research* 2022;11:889. doi: [10.12688/f1000research.122924.1](https://doi.org/10.12688/f1000research.122924.1)

80. Vivian J, Rao AA, Nothaft FA, Ketchum C, Armstrong J, Novak A, et al. Toil enables reproducible, open source, big biomedical data analyses. *Nature Biotechnology* 2017;35(4):314–316. doi: [10.1038/nbt.3772](https://doi.org/10.1038/nbt.3772)
81. ro-terms: Sapporo namespace [cited 2024 May 28].
<https://github.com/ResearchObject/ro-terms/tree/master/sapporo>
82. Suetake H, Fukusato T, Igarashi T, Ohta T. A workflow reproducibility scale for automatic validation of biological interpretation results. *GigaScience* 2023;12:giad031. doi: [10.1093/gigascience/giad031](https://doi.org/10.1093/gigascience/giad031)
83. Suetake H, Ohta TI, Tanjo T, Ishii M, Kinoshita BP, DrYak. sapporo-wes/sapporo-service: 1.5.1. Zenodo, 2023. doi: [10.5281/zenodo.10134452](https://doi.org/10.5281/zenodo.10134452)
84. Ohta T, Suetake H. Example of Workflow Run RO-Crate Output in Sapporo. Zenodo, 2023. doi: [10.5281/zenodo.10134581](https://doi.org/10.5281/zenodo.10134581)
85. Manubens-Gil D, Vegas-Regidor J, Prodhomme C, Mula-Valls O, Doblas-Reyes FJ. Seamless management of ensemble climate prediction experiments on HPC platforms. 2016 International Conference on High Performance Computing & Simulation (HPCS), 2016;895-900. doi: [10.1109/HPCSim.2016.7568429](https://doi.org/10.1109/HPCSim.2016.7568429)
86. Yoo AB, Jette MA, Grondona M. SLURM: Simple Linux Utility for Resource Management. *Job Scheduling Strategies for Parallel Processing (JSSPP 2003)*. Lecture Notes in Computer Science, 2003;2862. doi: [10.1007/10968987_3](https://doi.org/10.1007/10968987_3)
87. Feng H, Misra V, Rubenstein D. PBS: a unified priority-based scheduler. *Proceedings of the 2007 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, 2007;203–214. doi: [10.1145/1254882.1254906](https://doi.org/10.1145/1254882.1254906)
88. Bahra A. Managing work flows with ecFlow. *ECMWF Newsletter*, 2011;129:30–32. doi: [10.21957/nr843dob](https://doi.org/10.21957/nr843dob)
89. Oliver H, Shin M, Matthews D, Sanders O, Bartholomew S, Clark A, et al. Workflow Automation for Cycling Systems. *Computing in Science & Engineering* 2019;21(4):7–21. doi: [10.1109/MCSE.2019.2906593](https://doi.org/10.1109/MCSE.2019.2906593)
90. Beltrán Mora D, Castrillo M, Marciani MG, Kinoshita BP, Tenorio-Ku L, Gaya-Ávila A, et al. Autosubmit 4.0.100. Zenodo, 2023. doi: [10.5281/zenodo.10199020](https://doi.org/10.5281/zenodo.10199020)
91. Goble C, Cohen-Boulakia S, Soiland-Reyes S, Garijo D, Gil Y, Crusoe MR, et al. FAIR Computational Workflows. *Data Intelligence* 2020;2(1-2):108–121. doi: [10.1162/dint_a.00033](https://doi.org/10.1162/dint_a.00033)
92. Kinoshita BP. RO-Crate created using Autosubmit version 4.0.100 workflow running kinow/auto-mhm-test-domains. Zenodo, 2023. doi: [10.5281/zenodo.8144612](https://doi.org/10.5281/zenodo.8144612)
93. Samaniego L, Kumar R, Attinger S. Multiscale parameter regionalization of a grid-based hydrologic model at the mesoscale. *Water Resources Research*, 2010;46(5). doi: [10.1029/2008WR007327](https://doi.org/10.1029/2008WR007327)
94. Kumar R, Samaniego L, Attinger S. Implications of distributed hydrologic model parameterization on water fluxes at multiple scales and locations. *Water Resources Research* 2013;49(1):360–379. doi: [10.1029/2012WR012195](https://doi.org/10.1029/2012WR012195)
95. Leo S. Run of digital pathology tissue/tumor prediction workflow. Zenodo, 2023. doi: [10.5281/zenodo.7774351](https://doi.org/10.5281/zenodo.7774351)
96. The Galaxy Community. Galaxy. Version 23.1 Software Heritage Archive, 2023.
<https://identifiers.org/swh:1:rel:33ce0ce4f6e3d77d5c0af8cff24b2f68ba8d57e9>
97. Colonnelli I. StreamFlow run of digital pathology tissue/tumor prediction workflow. Zenodo, 2023. doi: [10.5281/zenodo.7911906](https://doi.org/10.5281/zenodo.7911906)

98. Del Rio M, Lianas L, Aspegren O, Busonera G, Versaci F, Zelic R, et al. AI Support for Accelerating Histopathological Slide Examinations of Prostate Cancer in Clinical Studies. Image Analysis and Processing. ICIAP 2022 Workshops. ICIAP 2022. Lecture Notes in Computer Science 2022;13373. doi: [10.1007/978-3-031-13321-3_48](https://doi.org/10.1007/978-3-031-13321-3_48)
99. CRS4 Digital Pathology Platform [cited 2024 May 27]. <https://github.com/crs4/DigitalPathologyPlatform>
100. RO-Crate profiles [cited 2024 July 1]. <https://www.researchobject.org/ro-crate/profiles.html#ro-crate-profiles>
101. MIRAX format [cited 2024 May 27]. <https://openslide.org/formats/mirax/>
102. Common Provenance Model RO-Crate profile [cited 2024 May 27]. <https://w3id.org/cpm/ro-crate>
103. Wittner R, Mascia C, Gallo M, Frexia F, Müller H, Plass M, et al. Lightweight Distributed Provenance Model for Complex Real-world Environments. Scientific Data 2022;9:503. doi: [10.1038/s41597-022-01537-6](https://doi.org/10.1038/s41597-022-01537-6)
104. Wittner R, Holub P, Mascia C, Frexia F, Müller H, Plass M. et al. Towards a Common Standard for Data and Specimen Provenance in Life Sciences. Learning Health Systems 2023;e10365. doi: [10.1002/lrh2.10365](https://doi.org/10.1002/lrh2.10365)
105. Wittner R, Gallo M, Leo S, Soiland-Reyes S. Packing provenance using CPM RO-Crate profile. Version 1.1. Zenodo, 2023. doi: [10.5281/zenodo.8095888](https://doi.org/10.5281/zenodo.8095888)
106. Wittner R, Soiland-Reyes S, Leo S, Meurisse M, Hermjakob H. BY-COVID D4.3 Provenance model for infectious diseases. Zenodo, 2024 doi: [10.5281/zenodo.10927253](https://doi.org/10.5281/zenodo.10927253)
107. The W3C SPARQL Working Group. SPARQL 1.1 Overview. W3C Recommendation 21 March 2013 [cited 2024 May 27]. <https://www.w3.org/TR/sparql11-overview/>
108. Ferreira da Silva R, Badia RM, Bala V, Bard D, Bremer PT, Buckley I, et al. Workflows Community Summit 2022: A Roadmap Revolution. arXiv:2304.00019, 2023. doi: [10.48550/arXiv.2304.00019](https://doi.org/10.48550/arXiv.2304.00019)
109. Rehm HL, Page AJH, Smith L, Adams JB, Alterovitz G, Babb LJ, et al. GA4GH: International policies and standards for data sharing across genomic research and healthcare. Cell Genomics 2021;1(2):100029. doi: [10.1016/j.xgen.2021.100029](https://doi.org/10.1016/j.xgen.2021.100029)
110. de Wit R. A Non-Intimidating Approach to Workflow Reproducibility in Bioinformatics: Adding Metadata to Research Objects through the Design and Evaluation of Use-Focused Extensions to CWLProv. Zenodo, 2022. doi: [10.5281/zenodo.7113250](https://doi.org/10.5281/zenodo.7113250)
111. de Wit R, Crusoe MR. Analysis of runcrate. Zenodo, 2024. doi: [10.5281/zenodo.12689424](https://doi.org/10.5281/zenodo.12689424)
112. Leo S, Crusoe MR, Rodríguez-Navas L, Sirvent R, Kanitz A, De Geest P, et al. Recording provenance of workflow runs with RO-Crate (RO-Crate and mapping). Zenodo, 2023. doi: [10.5281/zenodo.10368990](https://doi.org/10.5281/zenodo.10368990)
113. Leo S, Crusoe MR, Rodríguez-Navas L, Sirvent R, Kanitz A, De Geest P, et al. Recording provenance of workflow runs with RO-Crate (RO-Crate and mapping). HTML preview [cited 2024 May 27]. <https://w3id.org/ro/doi/10.5281/zenodo.10368989>
114. Soiland-Reyes S, Wheeler S. Five Safes RO-Crate profile. Version 0.4. TRE-FX Candidate Recommendation, 2023 [cited 2023 Dec 11]. <https://w3id.org/5s-crate/0.4>
115. Giles T, Soiland-Reyes S, Couldridge J, Wheeler S, Thomson B, Beggs J, et al. TRE-FX: Delivering a federated network of trusted research environments to enable safe data analytics. Zenodo, 2023. doi: [10.5281/zenodo.10055354](https://doi.org/10.5281/zenodo.10055354)

116. Desai T, Ritchie F, Welpton R. Five Safes: designing data access for research. Economics Working Paper Series, 2016;1601. <https://econpapers.repec.org/RePEc:uwe:wpaper:20161601>
117. Snowley K, Edwards L, Crosby B, Tatlow H. Integrating Our Community. Year 1. Health Data Research UK, 2023 (report) [cited 2023 Dec 11]. https://www.hdruk.ac.uk/wp-content/uploads/2023/10/Integrating-Our-Community_v1-Oct-2023-compressed.pdf
118. EOSC-ENTRUST: Creating a European network of TRUSTed research environments [cited 2024 May 27]. <https://eosc-entrust.eu/>
119. Mazumder R, Simonyan V (eds). IEEE P2791 BioCompute Working Group (BCOWG). IEEE Standard for Bioinformatics Analyses Generated by High-Throughput Sequencing (HTS) to Facilitate Communication. IEEE Std 2791-2020, 2020. doi: [10.1109/IEEESTD.2020.9094416](https://doi.org/10.1109/IEEESTD.2020.9094416)
120. Alterovitz G, Dean D, Goble C, Crusoe MR, Soiland-Reyes S, Bell A. Enabling Precision Medicine via standard communication of NGS provenance, analysis, and results. PLOS Biology 2018;16(12):e3000099. doi: [10.1371/journal.pbio.3000099](https://doi.org/10.1371/journal.pbio.3000099)
121. Stian Soiland-Reyes. Packaging BioCompute Objects using RO-Crate [cited 2024 May 27]. <https://biocompute-objects.github.io/bco-ro-crate/>
122. Soiland-Reyes S. Describing and packaging workflows using RO-Crate and BioCompute Objects. Zenodo, 2021. doi: [10.5281/zenodo.4633732](https://doi.org/10.5281/zenodo.4633732)