

# Recording provenance of workflow runs with RO-Crate

Simone Leo<sup>1\*</sup>, Michael R. Crusoe<sup>2,3,4</sup>, Laura Rodríguez-Navas<sup>5</sup>, Raül Sirvent<sup>5</sup>, Alexander Kanitz<sup>6,7</sup>, Paul De Geest<sup>8</sup>, Rudolf Wittner<sup>9,10,11</sup>, Luca Pireddu<sup>1</sup>, Daniel Garijo<sup>12</sup>, José M. Fernández<sup>5</sup>, Iacopo Colonnelli<sup>13</sup>, Matej Gallo<sup>9</sup>, Tazro Ohta<sup>14,15</sup>, Hirotaka Suetake<sup>16</sup>, Salvador Capella-Gutierrez<sup>5</sup>, Renske de Wit<sup>2</sup>, Bruno P. Kinoshita<sup>5</sup>, Stian Soiland-Reyes<sup>17,18</sup>

**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 Dept., 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, Inc., 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](mailto:simone.leo@crs4.it)

## 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 products (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 the PROV-O standard [6], an OWL [7] representation of PROV for provenance in the Web. PROV-O has been widely extended for workflows (D-PROV [8], ProvONE [9], OPMW [10], 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 can also be represented at different levels of abstraction depending on available computing resources: for instance, by the workflow execution becoming a single activity which produces results, by specifying the individual execution of each workflow step, or by going a step further and indicating how each step is divided into sub-processes when a workflow is deployed in a cluster.

Different workflow systems have adopted and extended PROV (and its PROV-O representation) to the workflow domain (WINGS [13, 14], VisTrails [15, 16]), in order to ease the burden of provenance collection from tool developers to workflow management systems (WMS) [17, 18].

D-PROV, PROV-ONE, OPMW-PROV, 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 like *wfprov* and *wfdesc* [19] go a step further by considering not only the link between plans and executions, but how to package the various artefacts as a Research Object (RO) [20] in order to ease portability while keeping the context of a digital experiment.

However, while these models address some workflow provenance representation issues, they have two main limitations: firstly, the extensions of PROV are not directly interoperable because of differences in granularity or different assumptions in their workflow representations; secondly, their support from WMS is typically one system per model. An early approach to unify and integrate workflow provenance traces across WMS was WEST (Workflow Ecosystems through STandards) [13], through the use of WINGS [14] 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 (<https://sheeba-samuel.github.io/REPRODUCE-ME/research/provbook.html>) and highlighting the need for provenance also where there is no workflow management system.

More recently, interoperability have been partially addressed by Common Workflow Language Prov (CWLProv) [22], which represents workflow enactments as ROs serialised according to the Big Data Bag (BDBag) approach [23]. The resulting format is a folder containing several data and metadata files [24], expanding on the RO Bundle approach of Taverna [25]. 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 specification [26] – a standard workflow specification adopted by at least a dozen different workflow systems (<https://www.commonwl.org/implementations/>). 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, as proposed by the *wfdesc* extension Roterms (<https://wf4ever.github.io/ro/2016-01-28/roterms>). In order 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 WMS. 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 proposed multiple levels of provenance [22, figure 2], from Level 0 (capturing workflow definition) to Level 3 (domain-specific annotations). In practice, the CWL reference implementation *cwltool* [27] and the corresponding CWLProv specification [24] records provenance details of all task executions together with the intermediate data and any nested workflows (CWLProv level 2), a granularity level that requires substantial support from the WMS. This approach is appropriate for workflow languages where the execution plan, including its distribution among the various tasks, is well known in advance (such as CWL). 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 [28]). This makes the implementation of CWLProv challenging, as shown by the fact that 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 (<https://w3id.org/bundle/context>) and stored in a separate manifest file, that includes the data identifier mapping to filenames. A project that uses serialised ROs similar to those used by CWLProv is Whole Tale [29], 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 [30] is a recent approach to packaging research data together with their metadata; it extends Schema.org [31], 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 [32] metadata file at the top level. The metadata file describes all entities stored in the RO-Crate along with their relationships; 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*. At the same time, the broad set of types and properties from Schema.org, complemented by a few additional terms from other vocabularies, make the RO-Crate model capable of expressing a wide range of contextual information that complements and enriches the core information specified by the profile. This may include, among others, the workflow authors and their affiliations, associated publications, licensing information, related software, etc. This is an approach used by WorkflowHub [33], a workflow system agnostic workflow registry which specifies a Workflow RO-Crate profile [34] 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 are the following:

- 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 [35], independent of the specific workflow language or execution system, and including support for reexecution.
- Implementations of the model in six workflow management systems and one conversion tool
- A mapping of our profiles against the W3C PROV-O Standard using the Simple Knowledge Organisation System (SKOS) [36]

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 WMS 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; we then illustrate implementations and usage examples; this is followed by a discussion and plans for future work.

## 2 The Workflow Run RO-Crate profiles

RO-Crate profiles are extensions of the base RO-Crate specification that describe how to represent the entities 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 [37].

The Workflow Run RO-Crate profiles are the main outcome of the activities of the Workflow Run RO-Crate Community (<https://www.researchobject.org/workflow-run-crate>), 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. In order to develop the Workflow-Run RO-Crate profiles, one of the first community efforts was to compile a list of requirements in the form of competency questions (<https://www.researchobject.org/workflow-run-crate/requirements>) to be addressed by the model. Each requirement was backed up 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. Many of the original issues are now closed, and the profiles have had four official releases on Zenodo.

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 execution 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 (<https://w3id.org/ro/wfrun/process>) to describe the execution of one or more tools that contribute to a computation; (2) Workflow Run Crate (<https://w3id.org/ro/wfrun/workflow>) to describe a computation orchestrated by a predefined workflow; (3) Provenance Run Crate (<https://w3id.org/ro/wfrun/provenance>) to describe a workflow computation including the internal details of individual step executions.

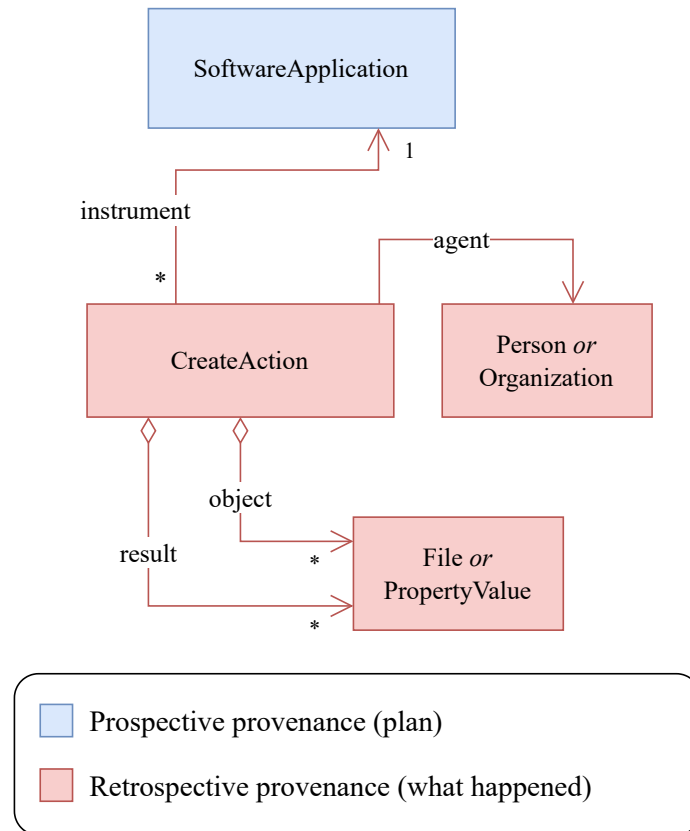
In the rest of this section we describe each of the above profiles in detail. We use italics to denote the types and properties describing entities and their relationships: these are defined in the RO-Crate JSON-LD context (<https://www.researchobject.org/ro-crate/1.1/context.jsonld>), which extends Schema.org with terms from the Bioschemas [38] ComputationalWorkflow profile (<https://bioschemas.org/profiles/ComputationalWorkflow/1.0-RELEASE>) and other vocabularies. More specifically, from Bioschemas we use the *ComputationalWorkflow* and *FormalParameter* types as well as the *input* and *output* properties. Note that these terms, though coming from Bioschemas, are not specific to the life sciences. We also developed a context extension through a dedicated “workflow-run” namespace (<https://w3id.org/ro/terms/workflow-run#>) to represent concepts that are not captured by terms in the RO-Crate context.

## 2.1 Process Run Crate

The Process Run Crate profile [39] contains specifications on describing 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. For instance, they could be executed manually by a human agent, one after the other as intermediate datasets become available, as shown in the process run crate (<https://w3id.org/ro/doi/10.5281/zenodo.6913045>) from [40]).

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

As an example, suppose a user called John Doe runs the `head` UNIX command to 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` command on `selection.txt`, storing the sorted output in a new file named `sorted.selection.txt`.

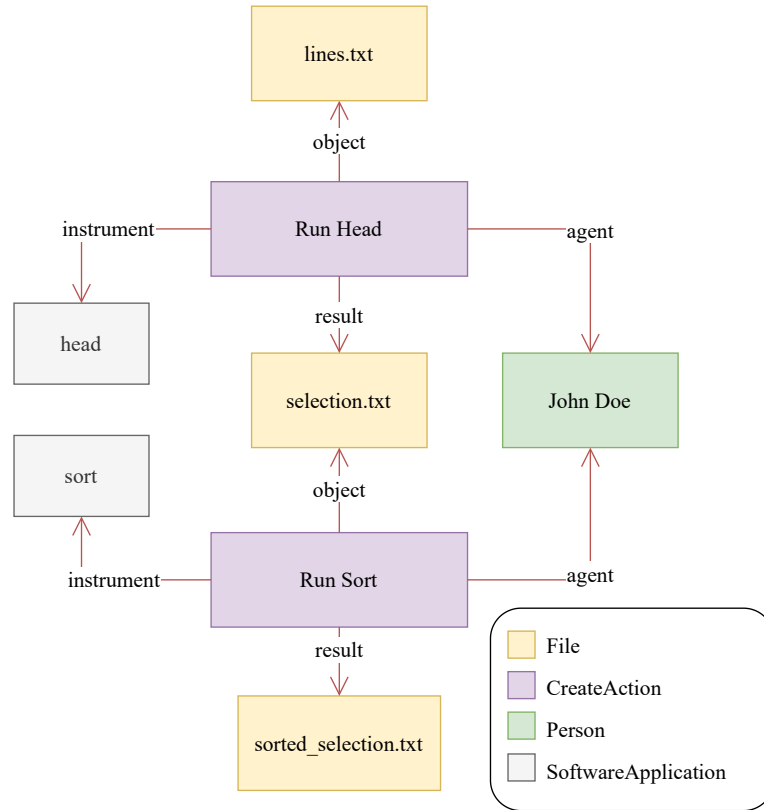


**Fig 1. UML class diagram for Process Run Crate.** The central entity is the *CreateAction*, which represents the execution of an application. It relates with the application itself via *instrument*, with the entity that executed it via *agent* and with its inputs and outputs via *object* and *result*, respectively. *File* is an RO-Crate alias for Schema.org’s *MediaObject*. Some inputs (and, less commonly, outputs), however, 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 *PropertyValue*. For simplicity, we left out the rest of the RO-Crate structure (e.g. the root *Dataset*). In this UML class notation diamond  $\diamond$  arrows indicate aggregation and regular arrows indicate references, \* indicates multiple instances, 1 means single instance.

Fig 2 contains a diagram of the two actions and their relationships to the other entities involved. 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.

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. One of the implications of this approach is that consumers need to code defensively, avoiding unwarranted





**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.

assumptions – e.g. by verifying that a value exists for an optional property before trying to retrieve it and use it.

## 2.2 Workflow Run Crate

The Workflow Run Crate profile [41] combines the Process Run Crate and WorkflowHub’s Workflow RO-Crate [34] profiles to describe the execution of “proper” computational workflows managed by a WMS. Such workflows are typically written in a special-purpose language, such as CWL or Snakemake [42], and run by one or more WMS (e.g., StreamFlow [43], Galaxy [44]). As in Process Run Crate, the execution is described by a *CreateAction* that links to the application via *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*, *SoftwareSourceCode* and *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 RO-Crate profile also contains recommendations on how to represent the workflow’s input and output parameters, based on the aforementioned Bioschemas [38] *ComputationalWorkflow* profile. All these elements are represented via the *FormalParameter* entity and are referenced from the main workflow via the *input*

and *output* properties. While the entities referenced from *object* and *result* in the *CreateAction* represent data entities and argument values that were actually used in the workflow execution, the ones referenced from *input* and *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 *exampleOfWork* property – the downloadable file is a realisation of the formal parameter definition. 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"}
}
```

The derivation of Workflow Run Crate from Workflow RO-Crate makes RO-Crates that conform to this profile compatible with the WorkflowHub workflow registry by also conforming to its Workflow RO-Crate profile. Thus, users of a WMS that implements this profile (or Provenance Run Crate, which inherits it) are able to register their workflows in WorkflowHub – together with an execution trace – by simply running them and uploading the resulting RO-Crates. Additionally, the inheritance mechanism allows to reuse the specifications already developed for Workflow RO-Crate, which form part of the guidelines on representing the prospective provenance.

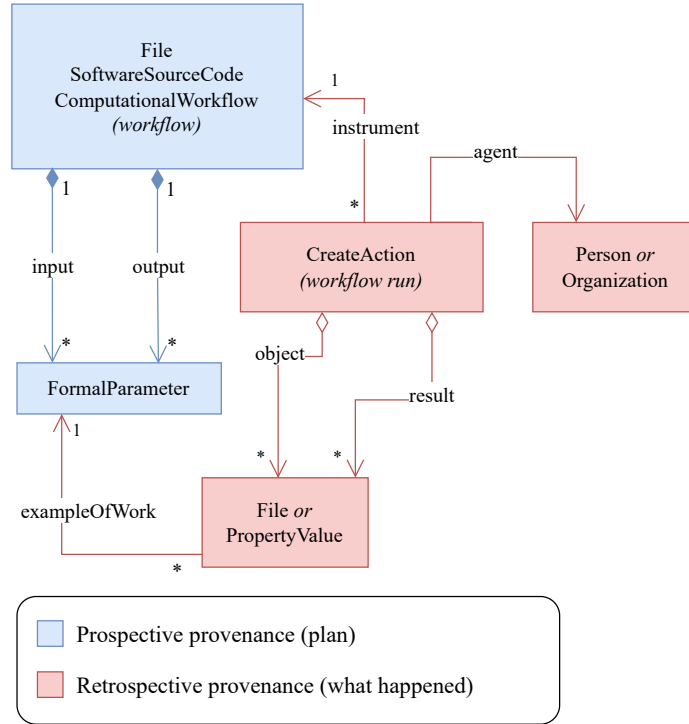
Fig 3 shows the entities used in Workflow Run Crate together with their relationships.

## 2.3 Provenance Run Crate

The Provenance Run Crate profile [45] 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 *CreateAction* instances that refer to the tool itself via *instrument* – analogously to its use in Process Run Crate. The workflow is required to refer to the tools it orchestrates through the *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 *HowToStep* i.e., “A step in the instructions for how to achieve a result” (<https://schema.org/HowToStep>). Steps point to the corresponding tools via the *workExample* property and are referenced from the workflow via the *step* property; the execution of a step is represented by a *ControlAction* pointing to the *HowToStep* via *instrument* and to the *CreateAction* instance(s) that represent the corresponding tool execution(s) via *object*. Note that a step execution does not coincide with a tool





**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 application is expected to be an entity with types *File*, *SoftwareSourceCode* and *ComputationalWorkflow*. Effectively, the entity belongs to all three types, and its properties are the union of the properties of the individual types. The filled diamond  $\blacklozenge$  indicates composition, empty diamond  $\diamond$  aggregation, and other arrows relations.

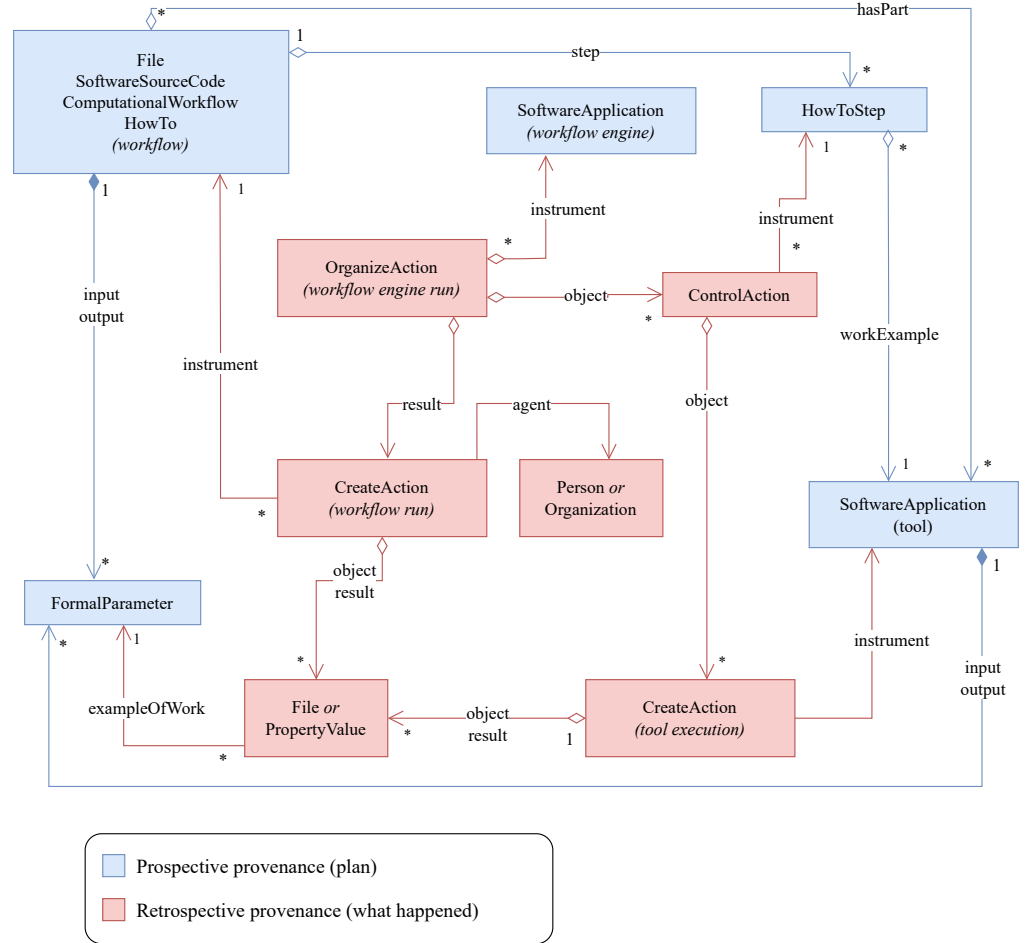
execution: an example where this distinction is apparent is when a step maps to multiple executions of the same tool over a list of inputs (e.g. the “scattering” feature in CWL).

An RO-Crate following this profile can also represent the execution of the WMS itself (e.g., cwltool) via *OrganizeAction*, pointing to a representation of the WMS via *instrument*, to the steps via *object* and to the workflow run via *result*. The *object* attribute of the *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 shows the various entities involved in the representation of a workflow run via Provenance Run Crate together with their relationships.

This profile also includes specifications on how to describe connections between parameters. Parameter connections – a fundamental feature of computational workflows – describe (i) how tools take as input the intermediate outputs generated by other tools and (ii) how workflow-level parameters are mapped to tool-level parameters. For instance, 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 connected to the input of the “head” tool wrapper, and the output of the latter is connected to the input of the “sort” tool wrapper.

A representation of parameter connections is particularly useful for traceability, since it allows 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 “workflow-run” context



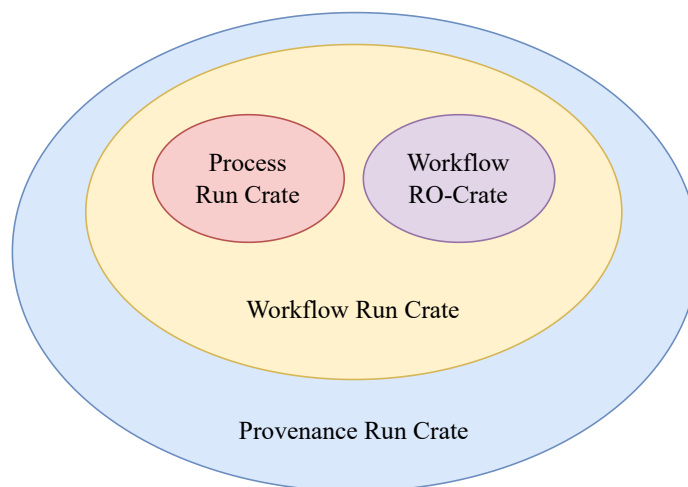
**Fig 4. UML class diagram for Provenance Run Crate.** In addition to the workflow run, this profile represents the execution of individual steps and their related tools.

extension (the <https://w3id.org/ro/terms/workflow-run#> namespace): a *ParameterConnection* type with *sourceParameter* and *targetParameter* attributes that respectively map to the source and target formal parameters, and a *connection* property to link from the relevant step or workflow to the *ParameterConnection* instances.

This profile is the most detailed of the three, and offers the highest level of granularity. Fig. 5 shows the relationship between the specifications of the profiles as a Venn diagram.

### 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. These tools have been developed in parallel by different teams, and independently from each other. RO-Crate has a strong ecosystem of tools [30], and these implementations have either re-used these or added their own approach to the standards.



**Fig 5. Venn diagram of the specifications for the various RO-Crate profiles.** 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.

### 3.1 Runcrate

Runcrate (<https://github.com/ResearchObject/runcrate>) [46] 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.

We also added a *run* subcommand to re-execute the computation described by an input Workflow Run Crate or Provenance Run Crate where CWL was used as a workflow language. It works by mapping the RO-Crate description of input parameters and their values (the workflow's *input* and the action'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 the characteristics of the workflow language (which should provide a clear formalism to map input items to their corresponding parameter slots) and from cooperation on the part of the workflow’s author, who can help considerably by containerizing the environment required by each step and providing the relevant annotations (if allowed by the workflow language).

## 3.2 Galaxy

The Galaxy project [44] 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 (<https://galaxyproject.org/toolshed/>); a web interface that allows the execution and definition of workflows using the integrated tools; a network of dedicated (public) Galaxy instances.

The export of workflow execution provenance data as Workflow Run Crates has been added in Galaxy’s 23.0 release. This feature provides a more interoperable alternative to the basic export of Galaxy workflow *invocations*: 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 is achieved by extracting provenance from Galaxy entities related to the workflow run, along with associated metadata, converting them to RO-Crate metadata using the ro-crate-py library [47]; by describing all files contained in the basic invocation export within the RO-crate metadata; and finally by making the Workflow Run Crate available for export to the user through Galaxy’s web interface and API [48].

We extract the prospective provenance contained in Galaxy’s YAML-based gxformat2 ([https://galaxyproject.github.io/gxformat2/v19\\_09.html](https://galaxyproject.github.io/gxformat2/v19_09.html)) workflow definition, which includes details of the analysis pipeline such as the graph of 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, parameter values used – is extracted from Galaxy’s data model (<https://docs.galaxyproject.org/en/master/lib/galaxy.model.html>), 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 exported Galaxy Workflow Run Crate is available on Zenodo [49].

In practice, a user would take the following steps to obtain a Workflow Run Crate from a Galaxy instance: (1) 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; (2) execute the workflow, providing the required inputs; (3) after the workflow has run successfully, the corresponding RO-Crate will be available for export from the Workflow Invocations list.

## 3.3 COMPSs

COMPSs [28] 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 making it efficient to exploit the resources available (either distributed or in a cluster). 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 (i.e. MPI, OpenMP), compilers (e.g. NUMBA), failure management, task grouping, and more.

Provenance recording for COMPSs workflows has been explored in previous work [50], where the Workflow RO-Crate profile was adopted in the implementation of a very lightweight approach to document provenance while avoiding the introduction of any significant run time performance overheads. However, because of the dynamic nature of COMPSs workflows, the Workflow Run Crate profile is better suited to represent them, since workflows are created when the application is executed and, thus, a prior static workflow definition does not exist before that moment. Due to this limitation, the workflow entity in the metadata file references the entry point application run by COMPSs, and formal parameters are not listed (note that listing them is not required by the profile) because inputs and outputs (both for each task and the whole workflow) are determined at runtime. COMPSs is able to export provenance data with a post-processing operation that can be triggered at any moment after the application has finished. The RO-Crate generation post-process uses information recorded by the runtime to detect and automatically add metadata of any input or output data assets used by the workflow.

Implementing Workflow Run Crate support in COMPSs required particular attention to the generation of a unique id for the *CreateAction* representing the workflow run, combining hostname and queuing system job id for supercomputer executions (as extra information added), and just using generated UUIDs for distributed environments, to add as much information as available from the run while ensuring the id is unique. In the *CreateAction*, the *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 *subjectOf* property of the *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 (<https://bsc.es/supportkc/docs/MareNostrum4/intro/>).

To showcase the COMPSs adoption of the Workflow Run Crate profile, we provide as an example the execution of the BackTrackBB [51] 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 [52] complies with the Workflow Run Crate profile, and includes the application source files, a diagram of the workflow's graph, application profiling and input and output files.

The implementation of provenance recording following Workflow Run Crate has been fully integrated in the COMPSs runtime, and is available since release 3.2 [53] (<https://github.com/bsc-wdc/compss/tree/3.2>).

### 3.4 StreamFlow

The StreamFlow framework [43] (<https://github.com/alpha-unito/streamflow>) is a container-native WMS based on the CWL standard. 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 [54], HPC batch scripts, and Helm charts [55]). 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 stores prospective and retrospective provenance data in a proprietary format into a persistent pluggable database (using sqlite3 as the default choice). After a CWL workflow execution completes, users can generate an RO-Crate through the `streamflow prov <workflow_name>` 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 [56].

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 (<https://github.com/inab/WfExS-backend>) 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 available software containers for reproducibility. Actual workflow execution is delegated to one of the supported workflow engines which matches with the workflow, right now either Nextflow or cwltool. 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 referred through resolvable identifiers, ideally public, permanent ones. Due to this, the orchestrator can consume workflows which are originally available in different kinds of locations, like git repositories, Software Heritage, or even RO-Crates from WorkflowHub.

WfExS-backend development milestones aim to reach FAIR workflow execution through the generation and consumption of RO-Crates following the latest Workflow Run Crate profiles, which have proven to be a mechanism suitable to semantically describe digital objects in a way that simplifies embedding key details involved in analysis reproducibility and replicability.

The orchestrator records details relevant to the prospective provenance when a workflow is prepared for execution, 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 metadata is represented 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. Thus, run details from all the executions are represented in the retrospective provenance RO-Crate. Support for Workflow Run RO-Crate is available since WfExS-backend



version 0.10.1 [77]. Future developments will also add support for embedding URLs of output results that have been deposited into a suitable repository (like Zenodo DOIs, for instance) as well as consuming previously produced RO-Crates.

An example of Workflow Run Crate generated by WfExS-backend from a Nextflow workflow run [58] is available from Zenodo [57].

### 3.6 Sapporo

Sapporo [59] is an implementation of the Workflow Execution Service (WES) API specification (<https://www.ga4gh.org/product/workflow-execution-service-wes/>). 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 [27], Toil [60], and StreamFlow [43]. 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 (<https://github.com/ResearchObject/ro-terms/tree/master/sapporo>).

Further, the Tonkatz companion command line software has integrated functionality to compare Run Crates generated by Sapporo to measure the reproducibility of the workflow outputs [61]. 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. Thus, 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 since release 1.5.1 [62]. An example of RO-Crate generated by Sapporo is available on Zenodo [63].

### 3.7 Autosubmit

Autosubmit [64] is an open source lightweight workflow manager and meta-scheduler created in 2011 for use in climate research to configure and run scientific experiments. It supports scheduling jobs via SSH to Slurm [65], PBS [66] and other remote batch servers used in HPC.

The “archive” feature was added in Autosubmit 3.1.0, released in 2015 (<https://earth.bsc.es/gitlab/es/autosubmit/-/tags/v3.1.0>). This 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 contains no structure, and users have no way to tell which is which from just looking at the ZIP file and its contents.

Recent releases of Autosubmit 4 include an updated YAML configuration management system that allows users to combine multiple YAML files into a single

unified configuration file. While this gave users flexibility, it also increased the complexity to track the configuration changes and to relate these to the provenance data. Another feature added in Autosubmit 4 is the option to use only the experiment manager features of Autosubmit, delegating the workflow execution to a different backend workflow engine, like ecFlow [67], Cylc [68], or a CWL runner.

In order to give users a more structured way to archive provenance, which includes the complete experiment configuration and the parameters used to generate the unified experiment configuration, and also to allow interoperability between workflow managers, the archive feature received a new flag in Autosubmit 4.0.100 [69] to generate Workflow Run RO-Crates.

The prospective provenance data is extracted from the Autosubmit experiment configuration. This includes the multiple YAML files, and 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 JSON-LD metadata file. Autosubmit uses the Workflow Run RO-Crate profile.

As one of the most recent implementations, much of the code added in Autosubmit 4 for RO-Crates was adapted from existing implementations like COMPSs and StreamFlow. `ro-crate-py` [47] was used for the heavy lifting work of creating the RO-Crate archive in Python, and adding information for the JSON-LD metadata.

The main challenges for adopting RO-Crate in Autosubmit were incorporating Autosubmit’s “Project” feature, and the lack of validation tools and of documentation and examples on how to use the standard with *coarse-grained* workflow management systems (as described in [70]) that do not track inputs and outputs, which is the case of Autosubmit – as well as the Cylc and ecFlow workflow engines.

A Project in Autosubmit is an abstract concept that has a type and a location, and is used to separate experiment configuration and template scripts and other auxiliary files. The type can be Git, Subversion, or Local. For each type the location represents the URL of a code repository, or a directory on a workstation or HPC file system used to copy the Project and its template scripts (written in Shell, R, or Python) and any other files (input data for a model, extra configuration files, binaries, etc.). The workflows in Autosubmit have tasks with dependencies to other tasks, and each of these tasks execute one of these template scripts. The RO-Crate file generated by Autosubmit includes only the project type and location, and not the complete Project. Therefore, users have the provenance of the Project, but only those with the correct permissions can access its constituent resources (many applications run with Autosubmit can not be publicly shared without consent).

Validation tools for RO-Crate archives are still under development, and while there is a community-based review process to help and guide new implementations, a tool that others can use as code is written will contribute to a more agile development.

After working with the RO-Crate community on these issues, 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. A pull request was created to `ro-crate-py` 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 missing information from the Autosubmit configuration model, like licence, authors, inputs, outputs, formal parameters, and more. And by modifying `ro-crate-py`, future implementers of RO-Crate that have a similar workflow configuration as Autosubmit should be able to re-use it, while also using COMPSs, StreamFlow, Autosubmit, and

other implementations as reference.

A workflow was created using an example Autosubmit Project [71] designed using UFZ’s mHM (mesoscale Hydrological Model) [72, 73]. This workflow was used to validate the RO-Crate produced by Autosubmit. This validation was performed by the Workflow Run RO-Crate community in a public GitHub repository (<https://github.com/ResearchObject/workflow-run-crate/>) and also using the aforementioned Runcrate.

### 3.8 Summary of implementations

Table 1 shows an overview of the different implementations presented in this section.

Table 1. Workflow Run Crate implementations

Impl.	Profile	Version URL/DOI	Example
runcrate	Provenance	[46]	[74]
Galaxy	Workflow	[75]	[49]
COMPSs	Workflow	[53]	[52]
Streamflow	Provenance	[56]	[76]
WfExS	Workflow	[77]	[57]
Sapporo	Workflow	[62]	[63]
Autosubmit	Workflow	[69]	[71]

Summary of each WRROC implementation, together with the profiles it implements, the latest software citation and an example crate of its application. Runcrate is a toolkit that converts CWLProv ROs to Provenance Run Crates, while the others are WMS.

## 4 Exemplary Use Cases

We illustrate Workflow Run RO-Crate on two exemplary use cases, which are similar in terms of application domain – machine learning-aided tumour detection in human prostate images – but 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, a combination of Process Run Crate and CPM RO-Crate is used to represent a sequence of computations linked to their corresponding CPM provenance information.

### 4.1 Provenance Run Crate for Digital Pathology

We present a use case that demonstrates the effectiveness of our most detailed profile Provenance Run Crate at recording 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 [78]. 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: inference of a low-resolution tissue mask to select areas for further processing; high-resolution tissue inference on areas identified in the previous step to refine borders; high-resolution cancer identification on areas identified in the first step. The two tissue inference steps run the same tool, but set different values for the parameter that controls the magnification level. The workflow is integrated in the CRS4 Digital Pathology Platform (<https://github.com/crs4/DigitalPathologyPlatform>), 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 same exemplary workflow in two different execution platforms. In the first case, the workflow was executed with the StreamFlow WMS, for which the Provenance Run Crate implementation is discussed in Section 3.4. In the second 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).

The RO-Crates obtained in the two cases [74, 76] are very similar to each other, differing only in a few details: for instance, [76] includes the StreamFlow configuration file and has separate files for the workflow and the two tools, while [74] has the workflow and the tools stored in a single file (CWL’s “packed” format). Apart from these minor differences, the description of the computation is essentially the same. 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 *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 entities and *PropertyValue* instances corresponding to the input and output values link to the corresponding parameter slots via the *exampleOfWork* property, providing information on the values taken by the parameters. The listing below (Fig 6) shows the output of the **runcrate report** command for the StreamFlow RO-Crate. For each action (workflow or tool run), the tool reports the associated instrument (workflow or tool), the starting and ending time and the list of inputs and outputs, with arrows pointing from the formal parameter to the corresponding actual value taken during the execution of the action.

The *exampleOfWork* link between input / output values and parameter slots is used by **runcrate run** to reconstruct the CWL input parameters document needed to rerun the computation. The *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, this particular workflow expects the input dataset to be in the MIRAX (<https://openslide.org/formats/mirax/>) format, where the “main” file taken as an input parameter by the processing application must be accompanied by a directory in the same location with the same name apart from the extension. The runcrate tool uses the *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 *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 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 best practices such as containerisation are followed.

This use case highlighted the need to add specifications on how to represent multi-file datasets [39, section Representing multi-file objects]. In the MIRAX format, in fact, the “main” file must be accompanied by a directory in the same location containing additional files with a specific structure. To represent this, we added

```

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

```

**Fig 6.** runcrate report command line output. This informal listing of relevant RO-Crate entities describe each step 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.

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 *Collection* entity linking to individual files and directories via the *hasPart* property, and referencing the main part (if any) via the *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 *Collection* entity that describes the whole dataset.

## 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. The RO-Crate makes use of Process Run Crate and CPM RO-Crate (<https://w3id.org/cpm/ro-crate>), an RO-Crate profile that supports the representation of entities described according to the Common Provenance Model (CPM) [79,80]. The CPM, an extension of the W3C PROV model [1] is a recently developed provenance model that enables the representation of distributed provenance. Distributed provenance 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. Individual provenance components are generated, stored, and managed individually by each organisation, and are linked together in a chain. 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, allowing to pack data, metadata, and CPM-based provenance information 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 pipeline consists of three main computational steps: a preprocessing step that splits input images into small patches and divides them into a training and a testing set; a training step that trains the model to recognise the presence of carcinoma cells in the images; an evaluation step that measures the accuracy of the trained model on the testing set. In addition to the 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 allowed us to describe all involved elements as a single aggregate, with entities and their relationships represented according to the RO-Crate model. The RO-Crate discussed here is available from Zenodo [81].

The CPM files complement the RO-Crate with internal details about the pipeline



execution, 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, it contains a representation of a checkpoint of the AI model after the second training iteration. The corresponding entity's attributes contain 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: (1) a biological sample acquired from a patient; (2) slices of the sample processed and put on glass slides; (3) 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 here 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 associated entities (inputs, outputs, code, etc.). The profiles are flexible, allowing to tailor the description to a broad variety of use cases, agnostic with respect to the WMS used and allow describing provenance traces at different levels of granularity. This facilitates developing implementations by multiple workflow systems (often with heterogeneous assumptions and requirements) – six of which have already been developed and are described in Section 3 – allowing to perform comparisons between runs across heterogeneous systems. For instance, the following SPARQL (<https://www.w3.org/TR/sparql11-overview/>) query returns all actions in a Workflow Run RO-Crate, together with their instruments and their starting and ending times:

```
PREFIX schema: <http://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 }
}
```

Additionally, 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.) rather than the trace alone. Being based on RO-Crate, the profiles and their implementations are part of a growing ecosystem of tools and services maintained by the RO-Crate community (<https://www.researchobject.org/ro-crate/in-use/>).

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 is mostly relevant in situations where the file is very large or cannot be shared publicly: the data entity’s identifier can be a URI that is accessible only through authentication, or resolvable only within the boundaries of the generating organisation.

The derivation of Workflow Run Crate from Workflow RO-Crate and, in turn, of Provenance Run Crate from Workflow Run Crate makes RO-Crates that conform to these profiles compatible with the WorkflowHub workflow registry, allowing workflow runs to be registered 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 Workflows Community Summit [82] identified as one of the current open challenges in the Scientific Workflows domain the ability to build FAIR into Workflow Management Systems, with the objective of achieving FAIR Computational Workflows. The profiles introduced in this article are able to help tackle this by introducing interoperable metadata among WMSs that captures the provenance of their corresponding workflow executions.

The Workflow Run RO-Crate profiles, the associated tooling, the implementations and the examples are developed by a community that runs regular virtual meetings (every two weeks at the time of writing) and coordinates on Slack and the RO-Crate mailing list. The WRROC community brings together members of the RO-Crate community [30], WMS users and developers, Workflow users and developers, GA4GH [83] Cloud developers and provenance model authors, and is open to anyone who is interested in the representation of workflow provenance. The inclusion of WMS developers and workflow users was key to keeping the specifications flexible, easy to implement and grounded on real use cases, while the diversity of the stakeholders allowed to keep a plurality of viewpoints while driving the model’s development forward.

One of the main benefits of this development process is that the profiles are already in use, with seven implementations (six WMS and one conversion tool) already available as described in section 3.

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

## 5.1 Evaluation of metadata coverage using `runcrate convert`

In order to assess the metadata coverage of `runcrate`, we performed a qualitative analysis of the tool’s *convert* mode, in which we evaluated how the generated RO-Crates preserve the metadata contained in the CWLProv ROs from which they are derived. For this analysis, we followed the same approach as for an earlier evaluation of CWLProv [84]. In that work, we identified and analysed three levels of representation: firstly, in RDF; secondly, in a structured, but CWL-specific document; and finally, in an unstructured, human readable format. From this earlier analysis, we concluded that the CWLProv RDF representation of the workflow runs lacked many provenance metadata that was included in CWL-specific documents, such as the packed workflow and input parameter file. For example, the CWLProv RDF only contained the name of each workflow step, without including the link to the underlying CommandLineTool or nested Workflow that was executed; information that could be extracted from the packed workflow.

In our analysis of `runcrate`, we compared the CWLProv RDF provenance graph with the RO-Crate metadata file. The results of the analysis are summarised in Table 2. The three dots (...) in the WRROC column indicate that the concept is supported in an RO-Crate using existing schema.org vocabulary (e.g.

<https://schema.org/softwareHelp>) but is not required or recommended by the WRROC profiles. Overall, most of the information contained in CWLProv RDF is transferred to the RO-Crate metadata. In addition, the representation of some 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 format of input files (D2), which was partially represented in CWLProv RDF, is fully represented in RO-Crate.

In conclusion, our analysis shows that runcrate preserves most provenance metadata previously shown to be relevant in realistic RO use case scenarios. The full results of the analysis can be found in [85].

**Table 2. Summarised results of our qualitative analysis of runcrate.**

Type	Subtype	Name	CWL	CWLProv	RO-Crate	WRROC
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 compared RO-Crates with the CWLProv ROs from which they were generated. The analysis was based on a provenance taxonomy reflecting relevant provenance metadata based on realistic use cases for ROs associated with a real-life bioinformatics workflow [84]. CWL-specific documents are: `packed.cwl` (the workflow), `primary-job.json` (the inputs file), and `primary-output.json` (the outputs file). Since `packed.cwl` is also included in RO-Crate, we only considered how the metadata was represented in `ro-crate-metadata.json`.

For completeness we also show the theoretical capability of the Provenance Run Crate profile (WRROC column) assuming all its MUST/SHOULD requirements are complete. The categories in the first three columns are explained in [84].

**Legend:** ● fully represented ○ partially represented · missing or unstructured representation ... optional (e.g. schema.org attribute)

From this analysis it is worth highlighting the gaps and potential for Workflow Run RO-Crates. Several areas have been flagged by this study as important aspects of workflow metadata, such as Data Access (D3), Software Documentation (SW2) and Workflow Requirements (WF3). Many such aspects require human annotation and cannot be provided by workflow engines alone, although they may be propagated from workflow and tool definitions. Some areas like Consumed Resources (EX2) require additional terms to be defined, and are part of future work.

## 5.2 Workflow Run RO-Crate and the W3C PROV standard

Our aim is to be 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 entities 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 online (<https://github.com/ResearchObject/workflow-run-crate-paper/tree/main/mapping>).

**Table 3. Mapping from Workflow Run RO-Crate to equivalent W3C PROV concepts** using SKOS [36]. For instance, *CreateAction* has **broader** match PROV’s *Activity*, meaning that *Activity* is more general.

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

## 5.3 Five Safes Workflow Run Crate

The *Five Safes RO-Crate* [86] profile has been developed to extend the Workflow Run RO-Crate profile for use in Trusted Research Environments (TRE) in order to handle sensitive health data in federated workflow execution across TREs in the UK [87] and following the Five Safes Framework [88]. 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 profile 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 [89] and in the European EOSC-ENTRUST project for Trusted Research Environments (<https://esciencelab.org.uk/projects/eosc-entrust/>).

## 5.4 Biocompute Object RO-Crate

IEEE 2791-2020 [90], colloquially *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 [91]. 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 (<https://biocompute-objects.github.io/bco-ro-crate/>) was developed as a profile to combine both standards [92]. 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 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 Crate profiles.

## 6 Conclusion and Future Work

In this work we presented Workflow Run RO-Crate, a collection of RO-Crate profiles to represent the provenance of the execution of computational workflows at different levels of granularity. We described each profile and their corresponding implementations, shown how they apply to real use cases and described the community behind their development process. Workflow Run RO-Crate has already been adopted by six WMS, including Galaxy, StreamFlow and COMPSs. The flexibility of our model eases its implementation in more systems, allowing interoperability between their workflow run descriptions.

Workflow Run RO-Crate is an ongoing project driven by an open community. A natural consequence of this is that the profiles are not static entities, but keep being updated to cater for new requirements and use cases. In-progress features are tracked in the GitHub repository issues section (<https://github.com/ResearchObject/workflow-run-crate/issues>) and are open to discussion for the community. New features under discussion include a representation of the execution environment and recording workflow resource usage. The runcrate

toolkit is planned to be expanded both to better support the current features and to include new ones that may arise.

Many of the presented implementations will also develop new features. For example, the Galaxy implementation will add 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 through the implementation of a profile; and further developing features to allow for user-guided import of RO-Crates as Galaxy datasets, histories and workflows.

Finally, we are currently exploring the cloud execution of Workflow Run RO-Crates. On the one hand, the Workflow Execution Service (WES) specification is used by the Global Alliance for Genomics and Health (GA4GH) [83] to enable WMS-agnostic interpretation of workflows and scheduling of task execution. On the other hand, the Task Execution Service (TES) specification enables the execution of individual, atomic, containerized 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 (1) 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 (2) 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.

## Supporting information

- Process Crate profile [39]
- Workflow Run Crate profile [41]
- Provenance Run Crate profile [45]
- Workflow Run RO-Crate Introduction [93] (from Galaxy Smörgåsbord 2023)
- Implementations and examples (Table 1)

## 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.

The authors acknowledge funding from: Sardinian Regional Government through the XData Project (S.L., L.P.); Spanish Government (contract PID2019-107255GB) (R.S.); MCIN/AEI/10.13039/501100011033 (CEX2021- 001148-S) (R.S.); Generalitat de Catalunya (contract 2021-SGR-00412) (R.S.); European High-Performance Computing Joint Undertaking (JU) (No 955558) (R.S.); EU NextGenerationEU/PRTR (project eFlows4HPC) (R.S.); EU Horizon research and innovation programme under Grant agreement No 101058129 (DT-GEO) (R.S.); ELIXIR Platform Task 2022-2023 funding for Task “Container Orchestration” (A.K.); Research Foundation - Flanders (FWO) for ELIXIR Belgium (I000323N and I002819N) (P.D.G.); Life Science Database Integration Project, NBDC of Japan Science and Technology Agency; JSPS KAKENHI (Grant Number 20J22439); European Commission Horizon 2020 825575 (European Joint



Programme on Rare Diseases; SC1-BHC-04-2018 Rare Disease European Joint Programme Cofund) (L.R.N., J.M.F., S.C.G.), 955558 (eFlows4HPC), 823830 (BioExcel-2), 824087 (EOSC-Life) (S.L., L.R.N., P.D.G., R.W., L.P., J.M.F., S.C.G., S.S.R.); Horizon Europe 101046203 (BY-COVID) (S.L., L.R.N., P.D.G., R.W., L.P., J.M.F., S.C.G., S.S.R.), 101057388 (EuroScienceGateway) (P.D.G., J.M.F., S.C.G., S.S.R.), 101057344 (FAIR-IMPACT) (D.G., S.S.R.); UK Research and Innovation (UKRI) under the UK government’s Horizon Europe funding guarantee 10038963 (EuroScienceGateway), 10038992 (FAIR-IMPACT).

H.S. is founder and CEO of the Software company Sator Inc., Tokyo.

## 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

**José M. Fernández** Data Curation, Software, Writing – Original Draft preparation

**Iacopo Colonnelli** Data Curation, Software, Writing – Original Draft preparation

**Matej Gallo** Data Curation, Software

**Tazro Ohta** Data Curation, Software, Writing – Original Draft preparation

**Hiroataka Suetake** Data Curation, Software, Writing – Original Draft preparation

**Salvador Capella-Gutierrez** Funding Acquisition, Resources, Supervision, Writing – Original Draft preparation

**Renske de Wit** Software, Writing – Original Draft preparation, Writing – Review & Editing

**Bruno de Paula Kinoshita** Data Curation, Software, Writing – Original Draft preparation, Writing – Review & Editing

**Stian Soiland-Reyes** Conceptualization, Formal Analysis, Funding Acquisition, Investigation, Methodology, Resources, Software, Supervision, Visualization, Writing – Original Draft preparation, Writing – Review & Editing

## 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. 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)
4. 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)
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://purl1.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. 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)
14. 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)
15. 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)
16. 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)

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. 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)
23. 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)
24. 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)
25. 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)
26. 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)
27. 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)
28. 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)
29. 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)
30. 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)
31. 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)
32. 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/>

33. 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)
34. 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>
35. 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)
36. 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/>
37. 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)
38. 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/>
39. Workflow Run RO-Crate working group. Process Run Crate specification. Version 0.4. Zenodo, 2023. doi: [10.5281/zenodo.10203944](https://doi.org/10.5281/zenodo.10203944)
40. Meurisse M, Estupiñán-Romero F, González-Galindo J, Martínez-Lizaga N, Royo-Sierra S, Saldner S, et al. Federated causal inference based on real-world observational data sources: application to a SARS-CoV-2 vaccine effectiveness assessment. *BMC Medical Research Methodology* 2023;23:248. doi: [10.1186/s12874-023-02068-3](https://doi.org/10.1186/s12874-023-02068-3)
41. Workflow Run RO-Crate working group. Workflow Run Crate specification. Version 0.4. Zenodo, 2023. doi: [10.5281/zenodo.10203971](https://doi.org/10.5281/zenodo.10203971)
42. 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)
43. 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)
44. 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)
45. Workflow Run RO-Crate working group. Provenance Run Crate specification. Version 0.4. Zenodo, 2023. doi: [10.5281/zenodo.10203978](https://doi.org/10.5281/zenodo.10203978)
46. 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)
47. De Geest P, Droesbeke 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)
48. 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)
49. De Geest P. Run of an example Galaxy collection workflow. Zenodo, 2023. doi: [10.5281/zenodo.7785861](https://doi.org/10.5281/zenodo.7785861)
50. 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)
51. 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)
  52. 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)
  53. 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)
  54. 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)
  55. 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)
  56. Colonnelli I, Cantalupo B, Aldinucci M, Saitta G, Mulone A (2023): StreamFlow. Version 0.2.0.dev10. Software Heritage Archive, 2023. <https://identifiers.org/swh:1:rev:b2014add57189900fa5a0a0403b7ae3a384df73b>
  57. Fernández González JM. RO-Crate from staged WfExS working directory 047b6dfc-3547-4e09-92f8-df7143038ff4 (overbridging templon). Zenodo, 2023. doi: [10.5281/zenodo.10091550](https://doi.org/10.5281/zenodo.10091550)
  58. Bouyssié D, Altiner P, Capella-Gutierrez S, Fernández JM, Hagemmeijer 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)
  59. 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)
  60. 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)
  61. 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)
  62. 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)
  63. 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)
  64. 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)
  65. 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)

66. 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)
67. Bahra A. Managing work flows with ecFlow. ECMWF Newsletter, 2011;129:30–32. doi: [10.21957/nr843dob](https://doi.org/10.21957/nr843dob)
68. 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)
69. Beltrán Mora D, Castrillo M, Marciani MG, Kinoshita BP, Tenorio-Ku L, Gaya-Àvila A, et al. Autosubmit 4.0.100. Zenodo, 2023. doi: [10.5281/zenodo.10199020](https://doi.org/10.5281/zenodo.10199020)
70. 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)
71. 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)
72. 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)
73. 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)
74. Leo S. Run of digital pathology tissue/tumor prediction workflow. Zenodo, 2023. doi: [10.5281/zenodo.7774351](https://doi.org/10.5281/zenodo.7774351)
75. The Galaxy Community. Galaxy. Version 23.1 Software Heritage Archive, 2023. <https://identifiers.org/swh:1:rel:33ce0ce4f6e3d77d5c0af8cff24b2f68ba8d57e9>
76. 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)
77. Fernández JM, Rodríguez-Navas L, Muñoz-Cívico A, Iborra P, Lea D. WfExS-backend. Version 0.10.1. Zenodo, 2023. doi: [10.5281/zenodo.10068956](https://doi.org/10.5281/zenodo.10068956)
78. 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)
79. 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)
80. 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)
81. 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)
82. 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)
83. 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)



84. 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)
85. De Wit R, Crusoe MR. Analysis of runcrate. Zenodo, 2023. doi: [10.5281/zenodo.10251812](https://doi.org/10.5281/zenodo.10251812)
86. 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>
87. 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)
88. 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>
89. 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.hdr.uk/wp-content/uploads/2023/10/Integrating-Our-Community\\_v1-Oct-2023-compressed.pdf](https://www.hdr.uk/wp-content/uploads/2023/10/Integrating-Our-Community_v1-Oct-2023-compressed.pdf)
90. 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)
91. 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)
92. 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)
93. Leo S. Workflow Run RO-Crate Introduction. Galaxy Training Materials, 2023 [cited 2023 Dec 11]. <https://training.galaxyproject.org/training-material/topics/fair/tutorials/ro-crate-workflow-run-ro-crate/tutorial.html>