


Evaluation of tools for describing, reproducing and reusing scientific workflows

Philipp Diercks¹Dennis Gläser ²Ontje Lünsdorf³Michael Selzer ⁴Bernd Flemisch ²Jörg F. Unger ¹

1. Department 7.7 Modeling and Simulation, Bundesanstalt für Materialforschung und -prüfung (BAM), Berlin.


2. Lehrstuhl für Wasser- und Umweltsystemmodellierung, University of Stuttgart, Stuttgart.

3. Institut für Vernetzte Energiesysteme, Deutsches Zentrum für Luft- und Raumfahrt, Oldenburg.

4. Institut für Angewandte Materialien-MM, Karlsruher Institut für Technologie, Karlsruhe.

**Date Received:**

2022-12-05

Licenses:This article is licensed under: **Keywords:**

FAIR, reproducibility, scientific workflows, tool comparison, workflow management

Data availability:

Data can be found here:

<https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements>**Software availability:**

Software can be found here:

<https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements>

Abstract. In the field of computational science and engineering, workflows often entail the application of various software, for instance, for simulation or pre- and postprocessing. Typically, these components have to be combined in arbitrarily complex workflows to address a specific research question. In order for peer researchers to understand, reproduce and (re)use the findings of a scientific publication, several challenges have to be addressed. For instance, the employed workflow has to be automated and information on all used software must be available for a reproduction of the results. Moreover, the results must be traceable and the workflow documented and readable to allow for external verification and greater trust. In this paper, existing workflow management systems (WfMSs) are discussed regarding their suitability for describing, reproducing and reusing scientific workflows. To this end, a set of general requirements for WfMSs were deduced from user stories that we deem relevant in the domain of computational science and engineering. On the basis of an exemplary workflow implementation, publicly hosted at GitHub (<https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements>), a selection of different WfMSs is compared with respect to these requirements, to support fellow scientists in identifying the WfMSs that best suit their requirements.

1 Introduction

2 With increasing volume, complexity and creation speed of scholarly data, humans rely more
3 and more on computational support in processing this data. The “FAIR guiding principles for
4 scientific data management and stewardship” [41] were introduced in order to improve the ability
5 of machines to automatically find and use that data. FAIR comprises the four foundational
6 principles “that all research objects should be Findable, Accessible, Interoperable and Reusable
7 (FAIR) both for machines and for people”. In giving abstract, high-level and domain-independent

8 guidelines, the authors answer the question of what constitutes good data management. However,
9 the implementation of these guidelines is still in its infancy, with many challenges not yet
10 identified and some of which may not have readily available solutions [31]. Furthermore, efforts
11 are made towards an Internet of FAIR Data and Services (IFDS) [17], which requires not only
12 the data, but also the tools and (compute) services to be FAIR.

13 Data processing is usually not a single task, but in general (and in particular for computational
14 simulations) relies on a chain of tools. Thus, to achieve transparency, adaptability and repro-
15 ducibility of (computational) research, the FAIR principles must be applied to all components
16 of the research process. This includes the tools (i. e. *any* research software) used to analyze the
17 data, but also the scientific workflow itself which describes how the various processes depend on
18 each other. In a community-driven effort, the FAIR principles are applied to research software
19 and are extended to its specific characteristics by the FAIR for Research Software Working
20 Group [9]. For a discussion of how the FAIR principles should apply to workflows and workflow
21 management systems (WfMSs) we refer to the work by Goble et. al. [20].

22 In addition, in recent years there has been a tremendous development of different tools (see
23 e. g. <https://github.com/pditommaso/awesome-pipeline>) that aid the definition and
24 automation of computational workflows. These WfMSs have great potential in supporting the
25 goal above which is further discussed in section 1.1.

26 In this paper, we would like to discuss how WfMSs can contribute to the transparency, adaptability
27 and reproducibility of computational research, which are aspects that ultimately increase the
28 credibility of research results. Based on the authors' experience, user stories that are relevant in
29 the domain of computational science and engineering are defined. These user stories are then
30 used to extract a set of general requirements for WfMSs. Several different tools are compared
31 with respect to these requirements to support fellow scientists in identifying the tools that best
32 suit their requirements. The list of tools selected for comparison is subjective and certainly not
33 complete. However, a GitHub repository [16] providing an implementation of an exemplary
34 workflow for all tools and a short documentation with a link to further information was created,
35 with the aim to continuously add more tools in the future. Furthermore, by demonstrating how
36 the different tools could be used, we would like to encourage people to use WfMSs in their daily
37 work.

38 1.1 Introduction to workflow management systems

39 In this paper, we use the term *process* to describe a computation, that is, the execution of a
40 program to produce output data from input data. A process can be arbitrarily complex, but
41 from the point of view of the workflow, it is a single, indivisible step. A *workflow* describes
42 how individual processes relate to each other. Software-driven scientific workflows are often
43 characterized by a complex interplay of various pieces of software executed in a particular order.
44 The output of one process may serve as input to a subsequent process, which requires them to
45 be executed sequentially with a proper mapping of outputs to inputs. Other computations are
46 independent of each other and can be executed in parallel. Thus, one of the main tasks of WfMSs
47 is the proper and efficient scheduling of the individual processes.

48 As shown in fig. 1, each process in the workflow, just as the workflow itself, takes some input to

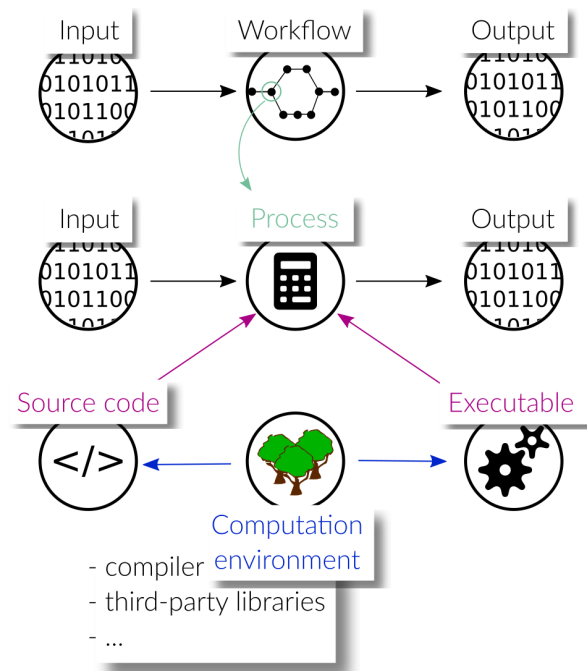


Figure 1: Schematic representation of software-driven scientific workflows.

49 produce output data. A more detailed discussion of the different levels of abstractions related to
 50 workflows can be found in Griem et al. [21]. The behavior of a process is primarily determined
 51 by the source code that describes it, but may also be influenced by the interpreters/compilers
 52 used for translation or the machines used for execution. Moreover, the source code of a process
 53 may carry dependencies to other software packages such that the behavior of a process possibly
 54 depends on their versions. We use the term *computation environment* to collect all those software
 55 dependencies, that is, interpreters and/or compilers as well as third-party libraries and packages
 56 that contribute to the computations carried out in a process. The exact version numbers of all
 57 involved packages are crucial, as the workflow may not work with newer or older packages, or,
 58 may produce different results.

59 As outlined in [30], WfMSs may be grouped into five classes. First, tools like *Galaxy* [1],
 60 *KNIME* [8], and *Pegasus* [14] provide a graphical user interface (GUI) to define scientific
 61 workflows. Thus, no programming skills are required and the WfMS is easily accessible for
 62 everybody. With the second group, workflows are defined using a set of classes and functions
 63 for generic programming languages (libraries and packages). This has the advantage that version
 64 control (e. g. using *Git* (<https://git-scm.com>)) can be employed on the workflow. In addition,
 65 the tool can be used without a graphical interface, e. g. in a server environment. Examples of
 66 prominent tools are *AiiDA* [23, 38], *doit* [34], *Balsam* [33], *FireWorks* [24] and *SciPipe* [28].
 67 Third, tools like *Nextflow* [15], *Snakemake* [27], *Bpipe* [32], *Guix Workflow Language* [44] and
 68 *Cluster Flow* [18] express the workflow using a domain specific language (DSL). A DSL is a
 69 language tailored to a specific problem. In this context, it offers declarations and statements to
 70 implement often occurring constructs in workflow definitions, which improves the readability

71 and reduces the amount of code. Moreover, the advantages of the second group also apply
72 for the third group. In contrast to the definition of the workflow in a programmatic way, the
73 fourth group comprises tools like *Popper* [25] and *Argo workflows* ([https://argoproj.g
74 ithub.io/argo-workflows/](https://argoproj.github.io/argo-workflows/)) which allow to specify the workflow in a purely declarative
75 way, by using configuration file formats like YAML [7]. In this case, the workflow specification
76 is concise and can be easily understood, but lacks expressiveness compared to the definition
77 of the workflow using programming languages. Fifth, there are system-independent workflow
78 specification languages like *CWL* [13] or *WDL* (<https://github.com/openwdl/wdl>). These
79 define a declarative language standard for describing workflows, which can then be executed by
80 a number of different engines like *Cromwell* [40], *Toil* [39], and *Tibanna* [29].

81 WfMSs can be used to create, execute and monitor workflows. They can help to achieve
82 reproducibility of research results by avoiding manual steps and automating the execution of
83 the individual processes in the correct order. More importantly, for a third person to reproduce
84 and reuse the workflow, it needs to be portable, that is, executable on any machine. Portability
85 can be supported by WfMSs with the integration of package management systems and container
86 technologies, which allow them to automatically re-instantiate the compute environment. Another
87 advantage of using WfMSs is the increase in transparency through a clear and readable workflow
88 specification. Moreover, after completion of the workflow, the tool can help to trace back
89 a computed value to its origin, by logging all inputs, outputs and possibly metadata of all
90 computations.

91 **2 User stories**

92 Starting from user stories that we consider representative for computational science and engi-
93 neering, a set of requirements is derived that serves as a basis for the comparison of different
94 WfMSs. In particular, a discussion on how the different tools implement these requirements is
95 provided.

96 Reproducibility, which is key to transparent research, is the main focus of the first user story
97 (see section 2.1). The second user story (see section 2.2) deals with research groups that develop
98 workflows in a joint effort where subgroups or individuals work on different components of the
99 workflow. Finally, the third user story focuses on computational research that involves generating
100 and processing large amounts of data, which poses special demands on how the workflow tools
101 organize the data that is created upon workflow execution (see section 2.3).

102 **2.1 Transparent and reproducible research paper**

103 *As a researcher, I want to share the code for my paper such that others are able to easily reproduce*
104 *my results.*

105 In this user story, the main objective is to guarantee the reproducibility of computational results
106 presented in scientific publications. Here, reproducibility means that a peer researcher is able to
107 rerun the workflow on some other machine while obtaining results that are in good agreement
108 with those reported in the publication. Mere reproduction could also be achieved without a
109 workflow tool, e. g. by providing a script that executes the required commands in the right

110 order, but this comes with a number of issues that may be solved with a standardized workflow
111 description.

112 First of all, reconstructing the logic behind the generation and processing of results directly from
113 script code is cumbersome and reduces the transparency of the research, especially for complex
114 workflows. Second, it is not straightforward for peer researchers to extract certain processes of a
115 workflow from a script and embed them into a different research project, hence the reusability
116 aspect is poorly addressed with this solution. Workflow descriptions may provide a remedy to
117 both of these issues, provided that each process in the workflow is defined as a unit with a clear
118 interface (see section 3.7).

119 While the workflow description helps peers to understand the details behind a research project,
120 it comes with an overhead on the side of the workflow creator, in particular when using a WfMS
121 for the first time. In the prevalent academic climate, we therefore think that an important aspect
122 of WfMSs is how easy they are to get started with (see section 3.9).

123 In the development phase, a workflow is typically run many times until its implementation is
124 satisfactory. With a scripted automation, the entire workflow is always executed, even if only one
125 process was changed since the last run. Since WfMSs have to know the dependencies between
126 processes, this opens up the possibility to identify and select only those parts of a workflow that
127 have to be rerun (see section 3.8). Besides this, the WfMS can display to the user which parts
128 are currently being executed, which ones have already been up-to-date, and which ones are still
129 to be picked (see section 3.2).

130 A general issue is that a workflow, or even each process in it, has a specific set of software- and
131 possibly hardware-requirements. This makes both reproducibility and reusability difficult to
132 achieve, especially over longer time scales, unless the computation environment in which the
133 original study was carried out is documented in a way that allows for a later re-instantiation.
134 The use of package managers that can export a given environment into a machine-readable
135 format from which they can then recreate that environment at a later time, may help to overcome
136 this issue. Another promising approach is to rely on container technologies. WfMSs have the
137 potential to automate the re-instantiation of a computation environment via integration of either
138 one of the above-mentioned technologies (see section 3.5). This makes it much easier for peers
139 to reproduce and/or reuse parts of a published workflow.

140 2.2 Joint research (group)

141 *As part of a research group, I want to be able to interconnect and reuse components of several*
142 *different workflows so that everyone may benefit from their colleagues' work.*

143 Similar to the previous user story, the output of such a workflow could be a scientific paper.
144 However, this user story explicitly considers interdisciplinary workflows in which the reusability
145 of individual components/modules is essential. Each process in the workflow may require a
146 different expertise and hence modularity and a common framework is necessary for an efficient
147 collaboration.

148 Many of the difficulties discussed in the previous user story are shared in a joint research project.
149 However, the collaborative effort in which the workflow description and those of its components

150 are developed promotes the importance of clear interfaces (see section 3.7) to ease communication
151 and an intuitive dependency handling mechanism (see section 3.5).

152 Another challenge here is that such workflows often consist of heterogeneous models of dif-
153 ferent complexity, such as large computations requiring high-performance computing (HPC),
154 preprocessing of experimental data or postprocessing analyses. Due to this heterogeneity, it may
155 be beneficial to outsource computationally demanding tasks to HPC systems, while executing
156 cheaper tasks locally (see section 3.1). Workflows with such computationally expensive tasks
157 can also strongly benefit from effective caching mechanisms and the reuse of cached results
158 wherever possible (see section 3.8).

159 Finally, support for a hierarchical embedding of sub-workflows (possibly published and versioned)
160 in another workflow is of great benefit as this allows for an easy integration of improve-
161 ments made in the sub-workflows by other developers (see section 3.6).

162 2.3 Complex hierarchical computations

163 *As a materials scientist, I want to be able to automate and manage complex workflows so I can*
164 *keep track of all associated data.*

165 Workflows in which screening or parameter sweeps are required typically involve running a large
166 number of simulations. Moreover, these workflows are often very complex with many levels of
167 dependencies between the individual tasks. Good data management that provides access to the
168 full provenance graph of all data can help to retain an overview over the large amounts of data
169 produced by such workflows (see section 3.4). For instance, the data management could be such
170 that desired information may be efficiently extracted via query mechanisms.

171 Due to the large amount of computationally demanding tasks in such workflows, it is helpful
172 if some computations can be outsourced to HPC systems (see section 3.1) with a clean way of
173 querying the current status during the typically long execution times (see section 3.2).

174 3 Specific requirements on workflow management systems

175 The user stories described above allow us to identify 11 requirements on WfMSs. They will be
176 described in the following and serve as evaluation criteria for the individual WfMSs discussed
177 in section 5.

178 3.1 Support for job scheduling system

179 As already mentioned, the main task of a WfMS is to automatically execute the processes of a
180 workflow in the correct order such that the dependencies between them are satisfied. However,
181 processes that do not depend on each other may be executed in parallel in order to speed up the
182 overall computation time. This requirement focuses on the ability of a workflow tool to distribute
183 the computations on available resources. Job scheduling systems like e. g. Slurm (also commonly
184 referred to as batch scheduling or batch systems) are often used to manage computations to be run
185 and their resource requirements (number of nodes, CPUs, memory, runtime, etc.). Therefore, it
186 is of great benefit if WfMSs support the integration of widely-used batch systems such that users

187 can specify and also observe the used resources alongside other computations that were submitted
 188 to their batch system in use. Besides this, this requirement captures the ability of a WfMS to
 189 outsource computations to a remote machine, e. g. a HPC cluster or cloud. For traditional HPC
 190 cluster systems it is usually necessary to transfer input and output data between the local system
 191 and the cluster system. This can be done using the secure shell protocol (SSH) and a WfMS may
 192 provide the automated transfer of a job's associated data. Ideally, the workflow can be executed
 193 anywhere without changing the workflow definition itself, but only the runtime arguments or a
 194 configuration file. The fulfillment of this requirement is evaluated by the following criteria:

- 195 The workflow system supports the execution of the workflow on the local system.
- 196 The workflow system supports the execution of the workflow on the local system via
 197 a batch system.
- 198 The workflow system supports the execution of the workflow via a batch system on
 199 the local or a remote system.

200 3.2 Monitoring



201 Depending on the application, the execution of scientific workflows can be very time-consuming.
 202 This can be caused by compute-intensive processes such as numerical simulations, or by a
 203 large number of short processes that are executed many times. In both cases, it can be very
 204 helpful to be able to query the state of the execution, that is, which processes have been finished,
 205 which processes are currently being executed, and which are still pending. A trivial way of such
 206 monitoring would be, for instance, when the workflow is started in a terminal which is kept
 207 open to inspect the output written by the workflow system and the running processes. However,
 208 ideally, the workflow system allows for submission of the workflow in the form of a process
 209 running in the background, while still providing means to monitor the state of the execution. For
 210 this requirement, two criteria are distinguished:

- 211 The only way to monitor the workflow is to watch the console output.
- 212 The workflow system provides a way to query the execution status at any time.

213 3.3 Graphical user interface



214 Independent of a particular execution of the workflow, the workflow system may provide
 215 facilities to visualize the graph of the workflow, indicating the mutual dependencies of the
 216 individual processes and the direction of the flow of data. One can think of this graph as the
 217 template for the data provenance graph. This visualization can help in conveying the logic
 218 behind a particular workflow, making it easier for other researchers to understand and possibly
 219 incorporate it into their own research. The latter requires that the workflow system is able
 220 to handle hierarchical workflows, that is, workflows that contain one or more sub-workflows
 221 as processes (see section 3.6). Beyond a mere visualization, a GUI may allow for visually
 222 connecting different workflows into a new one by means of drag & drop. We evaluate the
 223 features of a graphical user interface by means of the following three criteria:

- 224 The workflow system provides no means to visualize the workflow

- 225  The workflow system or third-party tools allow to visualize the workflow definition
- 226  The workflow system or third-party tools provide a GUI that enables users to graphi-
- 227 cally create workflows



228 3.4 Data provenance


229 The data provenance graph contains, for a particular execution of the workflow, which data and
 230 processes participated in the generation of a particular piece of data. Thus, this is closely related
 231 to the workflow itself, which can be thought of as a template for how that data generation should
 232 take place. However, a concrete realization of the workflow must contain information on the
 233 exact input data, parameters and intermediate results, possibly along with meta information on
 234 the person that executed the workflow, the involved software, the compute resources used and
 235 the time it took to finish. Collection of all relevant information, its storage in machine-readable
 236 formats and subsequent publication alongside the data can be very useful for future researchers
 237 in order to understand how exactly the data was produced. Ideally, the workflow system has the
 238 means to automatically collect this information upon workflow execution, which we evaluate
 239 using the following criteria:

- 240  The workflow system provides no means to export relevant information from a partic-
- 241 ular execution
- 242  The workflow system stores all results (also intermediate) together with provenance
- 243 metadata about how they were produced

244 3.5 Compute environment

245 In order to guarantee interoperability and reproducibility of scientific workflows, the work-
 246 flows need to be executable by others. Here, the re-instantiation of the compute environment
 247 (installation of libraries or source code) poses the main challenge. Therefore, it is of great
 248 use if the workflow tool is able to automatically deploy the software stack (on a per workflow
 249 or per process basis) by means of a package manager (e. g. conda <https://conda.io/>) or
 250 that running processes in a container (e. g. Docker <https://www.docker.com>, Apptainer
 251 <https://apptainer.org> (formerly Singularity)) is integrated in the tool. The automatic
 252 deployment of the software stack facilitates the execution of the workflow. However, it does not
 253 (always) enable reuse, that is, the associated software can be understood, modified, built upon
 254 or incorporated into other software [9]. For instance, if a container image is used, it is important
 255 that the container build recipe (e. g. Dockerfile) is provided. This increases the reusability as it
 256 documents how a productive environment, suitable to execute the given workflow or process,
 257 can be set up. The author of the workflow, however, is deemed to be responsible for the docu-
 258 mentation of the compute environment. For this requirement, we define the following evaluation
 259 criteria:


- 260  The automatic instantiation of the compute environment is not intended.
- 261  The workflow system allows the automatic instantiation of the compute environment
- 262 on a per workflow basis.


263  The workflow system allows the automatic instantiation of the compute environment
264 on a per process basis.

265 3.6 Hierarchical composition of workflows

266 A workflow consists of a mapping between a set of inputs (could be empty) and a set of outputs,
267 whereas in between a number of processes are performed. Connecting the output of one workflow
268 to the input of another workflow results in a new, longer workflow. This is particularly relevant
269 in situations where multiple people share a common set of procedures (e. g. common pre- and
270 postprocessing routines). In this case, copying the preprocessing workflow into another one is
271 certainly always possible, but does not allow to jointly perform modifications and work with
272 different versions. Moreover, a composition might also require to define separate compute
273 environments for each sub-workflow (e.g. using docker/singularity or conda). Executing all
274 sub-workflows in the same environment might not be possible because each sub-workflow might
275 use different tools or even the same tools but with different versions (e. g. python2 vs. python3).
276 This promotes the importance of supporting heterogeneous compute environments, which is
277 reflected in the evaluation criteria for this requirement:


278  The workflow system does not allow the composition of workflows.

279  The workflow system allows to embed a workflow into another one for a single
280 compute environment (homogeneous composition).


281  The workflow system allows to embed a workflow into another one for arbitrary many
282 (on a per process basis) compute environments (hierarchical composition).

283 3.7 Interfaces

284 In a traditional file-based pipeline, the output files produced by one process are used as inputs to
285 a subsequent process. However, it is often more convenient to pass non-file output (e. g. float or
286 integer values) directly from one process to another without the creation of intermediate files.
287 In this case, it is desirable that the workflow tool is able to check the validity of the data (e. g.
288 the correct data type) to be processed. Furthermore, this defines the interface for a process
289 more clearly and makes it easier for someone else to understand how to use, adapt or extend
290 the workflow/process. In contrast, in a file-based pipeline, this is usually not the case since a
291 dependency in form of a file does not give information about the type of data contained in that
292 file. We distinguish these different types of interfaces by the following criteria:

293  The workflow system is purely file-based and does not define interface formats.

294  The workflow system allows for passing file and non-file arguments between processes.

295  The workflow system allows for defining strongly-typed process interfaces, supporting
296 both file and non-file arguments.

297 3.8 Up-to-dateness

298 There are different areas for the application of workflows. On the one hand, people might use
299 a workflow to define a single piece of reproducible code that, when executed, always returns

300 the same result. Based on that, they might start a large quantity of different jobs and use the
 301 workflow system to perform this task. Another area of application is the constant development
 302 within the workflow (e.g. exchanging processes, varying parameters or even modifying the
 303 source code of a process) until a satisfactory result is obtained. The two scenarios require a
 304 slightly different behavior of the workflow system. In the first scenario, all runs should be kept
 305 in the data provenance graph with a documentation of how each result instance has been obtained
 306 (e.g. by always documenting the codes, parameters, and processes). If identical runs (identical
 307 inputs and processes should result in the same output) are detected, a recomputation should be
 308 avoided and the original output should be linked in the data provenance graph. The benefit of
 309 this behavior certainly depends on the ratio between the computation time for a single process
 310 compared to the overhead to query the data base.

311 However, when changing the processes (e.g. coding a new time integration scheme or a new
 312 constitutive model), the workflow system should rather behave like a build system (such as make)
 313 - only recomputing the steps that are changed or that depend on these changes. In particular for
 314 complex problems, this allows to work with complex dependencies without manually triggering
 315 computations and results in automatically recomputing only the relevant parts. An example is a
 316 paper with multiple figures where each is a result of complex simulations that in itself depend on
 317 a set of general modules developed in the paper. The “erroneous” runs are usually not interesting
 318 and should be overwritten.

319 How this is handled varies between the tools, yielding the following evaluation criteria:

320 **R** The complete workflow is always **R**ecomputed.

321 **L** A new entry in the data provenance graph is created which **L**inks the previous result
 322 (without the need to recompute already existing results).

323 **U** Only the parts are recreated (**U**psdated) that are not up-to-date. This usually reduces the
 324 overhead to store multiple instances of the workflow, but at the same time also prevents -
 325 without additional effort (e.g. when executing in different folders) computing multiple
 326 instances of the same workflow.

327 3.9 Ease of first use

328 Although this is not a requirement per-se, it is beneficial if the workflow system has an intuitive
 329 syntax/interface and little work is required for a new user to define a first workflow. Research
 330 applications typically have a high intrinsic complexity, and therefore, the complexity added by
 331 the workflow management should be as small as possible. Evaluation criteria:

332  difficult


333  intermediate


334  easy


335 3.10 Manually editable workflow definition

336 While it can be beneficial to create and edit workflows using a GUI (see section 3.3), it may be
 337 important that the resulting workflow description is given in a human-readable format. This

338 does not solely mean that the definition should be a text file, but also that the structure (e. g.
 339 indentation) and the naming are comprehensive. This facilitates version-controlling with git, in
 340 particular the code review process. Moreover, this does not force all users and/or developers to
 341 rely on the GUI. Evaluation criteria:

342  The workflow description is a binary file.

343  The workflow description is a text file but hard to interpret by humans.

344  The workflow description is a fully human-readable file format.

345 3.11 Platform for publishing and sharing workflows

346 The benefit of a workflow system is already significant when using it for individual research such
 347 as the development of an individual's paper or reproducing the paper someone else has written,
 348 when their data processing pipeline is fully reproducible, documented and published. However,
 349 the benefit can be even more increased if people are able to jointly work on (sub-)workflows
 350 together; particularly when a hierarchical workflow system is used. Even though workflows can
 351 easily be shared together with the work (e.g. in a repository), it might be beneficial to provide a
 352 platform that allows to publish documented workflows with a search and versioning functionality.
 353 This feature is not part of the requirement matrix to compare the different tools, but we consider
 354 a documentation of these platforms in the subsequent section as a good starting point for further
 355 research (exchange).

356 4 Simple use case

357 A simple exemplary use case was defined in order to analyze and evaluate the different workflow
 358 tools with respect to the requirements stated in section 3. This example is considered to be
 359 representative for many problems simulating physical processes in engineering science using
 360 numerical discretization techniques. It consists of six steps, as shown in fig. 2:

361 1. generation of a computational mesh (Gmsh)

362 2. mesh format conversion (MeshIO)

363 3. numerical simulation (FEniCS)

364 4. post-processing of the simulation results (ParaView)

365 5. preparation of macro definitions (Python)

366 6. compilation of a paper into a *.pdf* file using the simulation results (Tectonic)

367 The workflow starts from a given geometry on which the simulation should be carried out and
 368 generates a computational mesh in the first step using Gmsh [19]. Here, the user can specify the
 369 size of the computational domain by a float value `domain_size`. The resulting mesh file format
 370 is not supported by FEniCS [4], which is the software that we are using for the simulation carried
 371 out in the third step. Therefore, we convert the mesh file in the second step of the workflow from
 372 *.msh* to *.xdmf* using the python package MeshIO [35]. The simulation step yields result files in
 373 VTK file format [36] and returns the number of degrees of freedom used by the simulation as

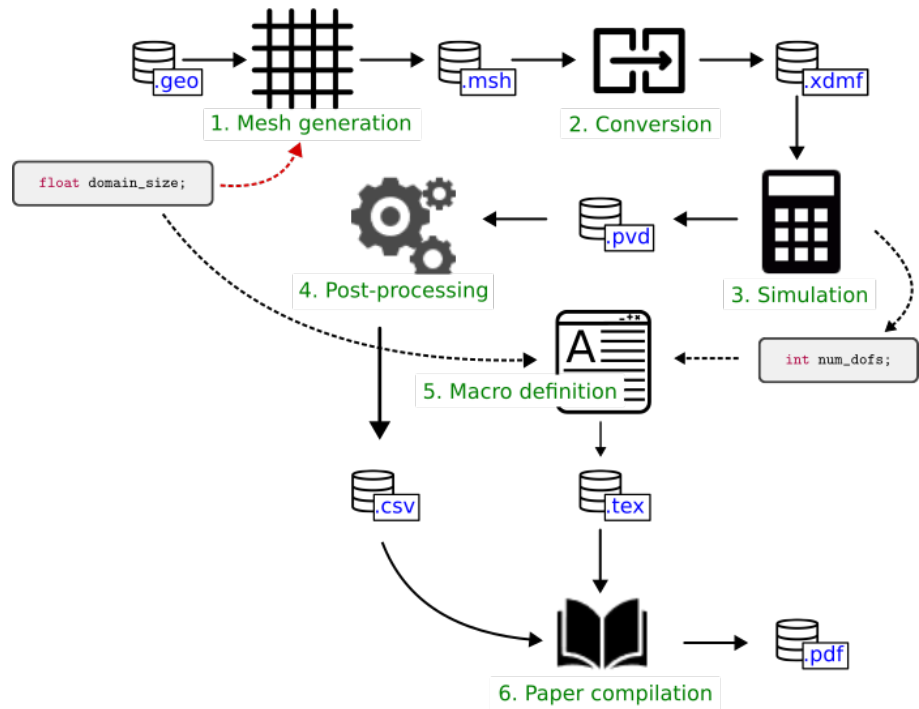


Figure 2: Task dependency graph of the simple use case. Mapping of input and output data is indicated with black arrows with solid lines. A dashed line refers to non-file input or output (parameters). Here, red color is used to distinguish user input from automatic data transfer.

374 an integer value `num_dofs`. The VTK files are further processed using the python application
 375 programming interface (API) of ParaView [2], which yields the data of a plot-over-line of the
 376 numerical solution across the domain in `.csv` file format. This data, together with the values for
 377 the domain size and the number of degrees of freedom, is inserted into the paper and compiled
 378 into a `.pdf` file using the \LaTeX engine Tectonic [42] in the final step of the workflow.

379 Most steps transfer data among each other via files, but we intentionally built in the transfer of
 380 the number of degrees of freedom as an integer value to check how well such a situation can
 381 be handled by the tools. Example implementations of the simple use case for various tools are
 382 available in a public repository [16].

383 5 Tool comparison

384 In this section, the selected WfMSs and their most important features are described and set in
 385 relation to the requirements defined in section 3. We note that to the best of our knowledge,
 386 existing add-on packages to the individual WfMSs are as well considered.

387 5.1 AiiDA

388 *AiiDA* [23, 38], the automated interactive infrastructure and database for computational science,
 389 is an open source Python infrastructure. With *AiiDA*, workflows are written in the Python
 390 programming language and managed and executed using the associated command line interface
 391 “*verdi*”.

392 *AiiDA* was designed for use cases that are more focused on running heavy simulation codes
393 on heterogeneous compute hardware. Therefore, one of the key features of *AiiDA* is the HPC
394 interface. It supports the execution of (sub-) workflows on any machine and most resource
395 managers are integrated. In case of remote computing resources, any data transfer, retrieval and
396 storing of the results in the database or status checking is handled by the *AiiDA* daemon process.
397 Another key feature is *AiiDA*'s workflow writing system which provides strongly typed interfaces
398 and allows for easy composition and reuse of workflows. Moreover, *AiiDA* automatically keeps
399 track of all inputs, outputs and metadata of all calculations, which can be exported in the form of
400 provenance graphs.

401 *AiiDA*'s workflow system enables to easily compose workflows, but *AiiDA* lacks in providing
402 the compute environment, such that the composition of heterogeneous workflows is challenging
403 since it requires the installation of software dependencies of the workflow on any machine that
404 should be used with *AiiDA*. The reason for this may be the challenges in using conda or containers
405 on HPC systems. On traditional HPC systems the preferred way of running software is to use
406 the provided module system to compile specific application code. The system may be isolated,
407 such that missing access to the internet prevents installing conda environments or downloading
408 container images. Moreover, successfully using container technology as an MPI-distributed
409 application across several nodes seems to be a technical challenge due to compatibility issues in
410 the MPI configuration and certain Infiniband drivers.

411 In addition to that, running external codes with *AiiDA* requires the implementation of an *AiiDA*
412 plugin which instructs *AiiDA* on how to run that code. This poses an additional overhead if the
413 application code changes frequently during development of the workflow. Also, in the special
414 case of FEniCS (see section 4), which can be used to solve partial differential equations and
415 therefore covers a wide spectrum of applications, it is very difficult to define a general plugin
416 interface which covers all models. We note that due to this use case which is rather different
417 from the use cases that *AiiDA* was designed for, the implementation of the simple use case
418 (see section 4) uses "aiida-shell" [22], an extension to the *AiiDA* core package which makes
419 running shell commands easy. While this is convenient to get a workflow running quickly, this
420 leads to an undefined process interface since this was the purpose of the aforementioned plugin
421 for an external code. Considering the points above, compared to the other tools, the learning
422 curve with *AiiDA* is fairly steep. In contrast to file-based workflow management systems, *AiiDA*
423 defines data types for any data that should be stored in the database. Consequently, non-file
424 based inputs are well defined, but this is not necessarily the case for file data.

425 In terms of the requirements defined in section 3, *AiiDA*'s strong points are execution, monitoring
426 and provenance. Due to the possibility to export provenance graphs, also level two of the
427 requirement "Graphical user interface" is reached. Lastly, caching can be enabled in *AiiDA*
428 to save computation time. Caching in *AiiDA* means, that the database will be searched for a
429 calculation of the same hash and if this is the case, the same outputs are reused.

430 5.2 Common Workflow Language

431 "*Common Workflow Language (CWL)* [5] is an open standard for describing how to run command
432 line tools and connect them to create workflows" (<https://www.commonwl.org/>). One

433 benefit of it being a standard is that workflows expressed in *CWL* do not have to be executed by
434 a particular workflow engine, but can be run by any engine that is able to parse the standard. In
435 fact, there exist a number of workflow engines that support *CWL* workflows, e. g. the reference
436 implementation *cwltool* (<https://github.com/common-workflow-language/cwltool>),
437 *Toil* [39] or *StreamFlow* [10].

438 *CWL* was designed with a focus on data analysis using command line programs. To create a
439 workflow, each of the command line programs is “wrapped” in a *CWL* description, defining what
440 inputs are needed, what outputs are produced and how to call the underlying program. Typically,
441 this step also reduces the possibly large number of options of the underlying command line tool
442 to a few options or inputs that are relevant for the particular task of the workflow. In a workflow,
443 the wrapped command line tools can be defined as individual processes, and the outputs of one
444 process can be mapped to the inputs of other processes. This information is enough for the
445 interpreter to build up the dependency graph, and processes that do not depend on each other
446 may be executed in parallel. A process can also be another workflow, thus, hierarchical workflow
447 composition is possible. Moreover, there exist workflow engines for *CWL* that support using job
448 managers like e. g. Slurm [45].

449 The *CWL* standard also provides means to specify the software requirements of a process. For
450 instance, one can provide the URL of a docker image or docker file to be used for the execution
451 of a process. In case of the latter, the image is automatically built from the provided docker file,
452 which itself contains the information on all required software dependencies. Besides this, the
453 *CWL* standard contains language features that allow listing software dependencies directly in
454 the description of a workflow or process, and workflow engines may automatically make these
455 software packages available upon execution. As one example, the current release of *cwltool*
456 supports the definition of software requirements in the form of e. g. *Conda* packages that are then
457 automatically installed when the workflow is run (see e. g. our implementation and the respective
458 pipelines at [16]).

459 In contrast to workflow engines that operate within a particular programming language, the
460 transfer of data from one process to another cannot occur directly via memory with *CWL*. For
461 instance, if the result of a process is an integer value, this value has to be read from a file produced
462 by the process, or, from its console output. However, this does not have to be done in a separate
463 process or by again wrapping the command line tool inside some script, since *CWL* supports the
464 definition of inline JavaScript code that is executed by the interpreter. This allows retrieving, for
465 instance, integer or floating point return values from a process with a small piece of code.

466 *CWL* requires the types of all inputs and outputs to be specified, which has the benefit that the
467 interpreter can do type checks before the execution of the workflow. A variety of primitive
468 types, as well as arrays, files or directories, are available. Files can refer to local as well as
469 online resources, and in the case of the latter, resources are automatically fetched and used upon
470 workflow execution.

471 There exist a variety of tools built around the *CWL* standard, such as the Rabix Composer (<https://rabix.io/>) for visualizing and composing workflows in a GUI. Besides that and as mentioned
472 before, there are several workflow engines that support *CWL* and some of which provide extra
473 features. For instance, *cwltool* allows for tracking provenance information of individual workflow
474

475 runs. However, to the best of our knowledge, there exists no tool that automatically checks which
476 results are up-to-date and do not have to be reproduced (see section 3.8).

477 The CWL standard allows to specify the *format* of an input or output file by means of an *IRI*
478 (Internationalized Resource Identifier) that points to online-available resource where the file
479 format is defined. For processes whose output files are passed to the inputs of subsequent jobs,
480 the workflow engine can use this information to check if the formats match. To the best of our
481 knowledge, *cwltool* does so by verifying that the *IRIs* are identical, or performs further reasoning
482 in case the *IRIs* point to classes in ontologies (see, for instance, the class for the JSON file format
483 in the EDAM ontology at edamontology.org/format_3464). Such reasoning can make use of
484 defined relationships between classes of the ontology to determine file format compatibility. For
485 more information on file format specifications in CWL see [commonwl.org/user_guide/topics/file-](https://commonwl.org/user_guide/topics/file-formats.html)
486 [formats.html](https://commonwl.org/user_guide/topics/file-formats.html).

487 5.3 *doit*

488 “*doit* comes from the idea of bringing the power of build-tools to execute any kind of task” [34].
489 The automation tool *doit* is written in the Python programming language. In contrast to systems
490 which offer a GUI, knowledge of the programming language is required. However, it is not
491 required to learn an additional API since task metadata is returned as a Python dictionary.
492 Therefore, we consider this as very easy to get started quickly.

493 With *doit*, any shell command available on the system or python code can be executed. This
494 also includes the execution of processes on a remote machine, although all necessary steps (e. g.
495 connecting to the remote via SSH) need to be defined by the user. In general, such behavior
496 as described in section 3.1 is possible, but it is not a built-in feature of *doit*. Also, *doit* does
497 not intend to provide the compute environment. Therefore, while in general the composition of
498 workflows (see section 3.6) is easily possible via python imports, this only works for a single
499 environment. The status of the execution can be monitored via the console. Here, *doit* will skip
500 the execution of processes which are up-to-date and would produce the same result of a previous
501 execution. To determine the correct order in which processes should be executed, *doit* also
502 creates a directed acyclic graph (DAG) which could be used to visualize dependencies between
503 processes using “*doit-graph*” (<https://github.com/pydoit/doit-graph>), an extension to
504 *doit*. For each run (specific instance of the workflow), *doit* will save the results of each process
505 in a database. However, the tool does not provide control over what is stored in the database.
506 On the one hand, *doit* allows to pass results of one process as input to another process directly,
507 without creating intermediate files, so it is not purely file-based. On the other hand, the interface
508 for non-file based inputs does not define the data type.

509 5.4 Guix Workflow Language

510 The *Guix Workflow Language (GWL)* [44] is an extension to the open source package manager
511 GNU Guix [12]. *GWL* leverages several features from Guix, chief among these is the compute
512 environment management. Like Guix, *GWL* only supports GNU/Linux systems.

513 *GWL* can automatically construct an execution graph from the workflow process input/output

514 dependencies but also allows a manual specification. Support for HPC schedulers via DRMAA¹
515 is also available.

516 *GWL* doesn't provide a graphical user interface, interactions are carried out using a command-line
517 interface in a text terminal. Monitoring is also only available in the form of simple terminal
518 output.

519 There is support to generate a GraphViz (see e. g. <https://graphviz.org>) description of the
520 workflow, which allows basic visualization of a workflow. Although not conveniently exposed²,
521 *GWL* has a noteworthy unique feature inherited from Guix: precise software provenance tracking.
522 Guix contains complete build instructions for every package (including their history through git),
523 which enables accounting of source code and the build process, like for example compile options,
524 of all tools used in the workflow. Integrity of this information is ensured through cryptographic
525 hash functions. This information can be used to construct data provenance graphs with a high
526 level of integrity (basically all userspace code of the compute environment can be accounted
527 for [11]).

528 *GWL* uses Guix to setup compute environments for workflow processes. Each process is
529 executed in an isolated³ compute environment in which only specified software packages are
530 available. This approach minimizes (accidental) side-effects from system software packages
531 and improves workflow reproducibility. Interoperability also benefits from this approach, since
532 a Guix installation is the only requirement to execute a workflow on another machine. As Guix
533 provides build instructions for all software packages, it should be easily possible to recreate
534 compute environments in the future, even if the originally compiled binaries have been deprecated
535 in the meanwhile (see [3] for a discussion about long-term reproducibility).

536 Composition of workflows is possible, workflows can be imported into other workflows. Com-
537 position happens either by extracting individual processes (repurposing them in a new workflow)
538 or by appending new processes onto the existing workflow processes.

539 *GWL* relies exclusively on files as interface to workflow processes. There's no support to
540 exchange data on other channels, as workflow processes are executed in isolated environments.

541 Like other workflow tools, *GWL* caches the result of a workflow process using the hash of its
542 input data. If a cached result for the input hash value exists, the workflow processes execution is
543 skipped.

544 *GWL* is written in the Scheme [37] implementation GNU Guile [43], but in addition to Scheme,
545 workflows can also be defined in wisp [6], a variant of Scheme with significant whitespace. wisp
546 syntax thus resembles Python, which is expected to flatten the learning curve a bit for scientific
547 audience. However, error messages are very hard to read without any background in Scheme.
548 On first use, *GWL* will be very difficult in general. This problem is acknowledged by the *GWL*
549 authors and might be subject to improvements in the future.

1. Distributed Resource Management Application API <https://www.drmaa.org>

2. *GWL* doesn't provide a command to export provenance graphs in any way, instead Guix needs to be queried for build instruction, dependency graphs and similar provenance information of a workflows software packages

3. By default, lightweight isolation is setup by limiting the PATH environment variable to the compute environment. Stronger isolation via Linux containers is also optionally available.

550 As both *wisp* and Scheme code is almost free of syntactic noise in general, workflows are almost
551 self-describing and easily human-readable.

552 In summary, *GWL* provides a very interesting and sound set of features especially for repro-
553 ducibility and interoperability. These features come at the cost of a Guix installation, which
554 requires administrator privileges. The workflow language is concise and expressive, but error
555 messages are hard to read. At the current stage, *GWL* can only be recommended to experienced
556 scheme programmers or to specialists with high requirements on software reproducibility and
557 integrity.

558 5.5 Nextflow and Snakemake

559 With *Nextflow* [15] and *Snakemake* [30], the workflow is defined using a DSL which is an
560 extension to a generic programming language (Groovy for *Nextflow* and Python for *Snakemake*).
561 Moreover, *Nextflow* and *Snakemake* also allow to use the underlying programming language
562 to generate metadata programmatically. Thus, authoring scientific workflows with *Nextflow* or
563 *Snakemake* is very easy.

564 The process to be executed is usually a shell command or an external script. The integration
565 with various scripting languages is an import feature of *Snakemake* as well as *Nextflow*, which
566 encourages readable modular code for downstream plotting and summary tasks. Also boilerplate
567 code for command line interfaces (CLIs) in external scripts can be avoided. Another feature of
568 *Snakemake* is the integration of Jupyter notebooks, which can be used to interactively develop
569 components of the workflow.

570 Both tools implement a CLI to manage and run workflows. By default, the status of the execution
571 is monitored via the console. With *Nextflow*, it is possible to monitor the status of the execution
572 via a weblog. *Snakemake* supports an external server to monitor the progress of submitted
573 workflows.

574 With regard to the execution of the workflow (section 3.1), the user can easily run the workflow
575 on the local machine and the submission via a resource manager (e. g. Slurm, Torque) is integrated.
576 Therefore, individual process resources can be easily defined with these tools if the workflow is
577 submitted on a system where a resource manager is installed, i. e. on a traditional HPC cluster
578 system. Despite this, only level two of the defined criteria is met for *Nextflow*, since the execution
579 of the workflow on a remote machine and the accompanied transfer of data is not handled by the
580 tool. For *Snakemake*, if the CLI option “default-remote-provider” is used, all input and output
581 files are automatically down- and uploaded to the defined remote storage, such that no workflow
582 modification is necessary.

583 The requirement “up-to-dateness” is handled differently by *Nextflow* and *Snakemake*. By default,
584 *Nextflow* recomputes the complete workflow, but with a single command-line option existing
585 results are retrieved from the cache and linked such that a re-execution is not required. In
586 this case, *Nextflow* allows storing multiple instances of the same workflow upon variation of a
587 configuration parameter. *Snakemake* will behave like a build tool in this context and skip the
588 re-execution of processes whose targets already exist and update any process whose dependencies
589 have changed.

590 A strong point of *Nextflow* and *Snakemake* is the integration of the conda package management
 591 system and container technologies like docker. For example, the compute environment can be
 592 defined for each process based on a conda environment specification file or a certain docker
 593 image. Upon execution of the workflow, the specified compute environment is re-instantiated
 594 automatically by the workflow tool, making it very easy to reproduce results of or built upon
 595 existing workflows. Furthermore, since the tool is able to deploy the software stack on a per
 596 process basis, the composition of hierarchical workflows as outlined in section 3.6 is possible.

597 Similar to *doit*, both tools do not provide a GUI to graphically create and modify workflows.
 598 However, a visualization of the workflow, i. e. a dependency graph of the processes, can be
 599 exported. Moreover, it is possible to export extensive reports detailing the provenance of the
 600 generated data.

601 *Nextflow* and *Snakemake* can also be regarded as file-based workflow management systems.
 602 Therefore, interface formats, i. e. class structures or types of the parameters passed from one
 603 process to the subsequent one, are not clearly defined.

604 5.6 Evaluation matrix

605 The evaluation of the workflow tools provided in section 5 in terms of the requirements described
 606 in section 3 on the example of the workflow outlined in section 4 yields the evaluation matrix
 depicted in table 1.

Table 1: Evaluation of the workflow tools.

Requirement	Workflow tool					
	<i>AiiDA</i>	<i>CWL</i>	<i>doit</i>	<i>GWL</i>	<i>Nextflow</i>	<i>Snake- make</i>
Job scheduling system	●●●	●●○	●○○	●●○	●●○	●●●
Monitoring	●●	●●	●○	●○	●○	●○
Graphical user interface	●●○	●●●	●●○	●○○	●●○	●●○
Provenance	●●●	●●○	●○○	●○○	●●○	●●○
Compute environment	●○○	●●●	●○○	●●●	●●●	●●●
Composition	●●○	●●●	●●○	●●●	●●●	●●●
Process interfaces	●●○	●●●	●○○	●○○	●○○	●○○
Up-to-dateness	L	R	U	U	L	U
Ease-of-first-use	●○○	●●○	●●●	●○○	●●●	●●●
Manually editable	●●●	●●●	●●●	●●●	●●●	●●●

607

608 6 Summary

609 In this work, six different WfMSs (*AiiDA*, *CWL*, *doit*, *GWL*, *Nextflow* and *Snakemake*) are studied.
 610 Their performance is evaluated based on a set of requirements derived from three typical user
 611 stories in the field of computational science and engineering. On the one hand, the user stories
 612 are focusing on facilitating the development process, and on the other hand on the possibility of
 613 reusing and reproducing results obtained using research software. The choice for one WfMS or
 614 the other is strongly subjective and depends on the particular application and the preferences of

615 its developers. The overview given in table 1 together with the assessments in section 5 may
616 only serve as a basis for an individual decision making.

617 For researchers that want to start using a WfMS, an important factor is how easy it is to get a first
618 workflow running. For projects that are written in Python, a natural choice may be *doit*, which
619 operates in Python and is easy to use for anyone familiar with the language. Another benefit
620 of this system is that one can use Python functions as processes, making it possible to easily
621 transfer data from one process to the other via memory without the need to write and read to
622 disk. In order to make a workflow portable, developers have to provide additional resources that
623 allow users to prepare their environment such that all software dependencies are met, prior to
624 the workflow execution.

625 To create portable workflows more easily, convenient tools are *Nextflow* or *Snakemake*, where
626 one can specify the compute environment in terms of a conda environment file or a container
627 image on a per-process basis. They require to learn a new domain-specific language, however,
628 our assessment is that it is easy to get started as only little syntax has to be learned in order to get
629 a first workflow running.

630 The strengths of *AiiDA* are the native support for distributing the workload on different (registered)
631 machines, the comprehensive provenance tracking, and also the possibility to transfer data among
632 processes without the creation of intermediate files.

633 *CWL* has the benefit of being a language standard rather than a specific tool maintained by a
634 dedicated group of developers. This has led to a variety of tooling developed by the community
635 as e. g. editors for visualizing and modifying workflows with a GUI. Moreover, the workflow
636 description states the version of the standard in which it is written, such that any interpreter
637 supporting this standard should execute it properly, which reduces the problem of version pinning
638 on the level of the workflow interpreter.

639 Especially for larger workflows composed of processes that are still under development, and
640 are thus changing over time, it may be useful to rely on tools that allow to define the process
641 interfaces by means of strongly-typed arguments. This can help to detect errors early on, e. g. by
642 static type checkers. *CWL* and *AiiDA* support the definition of strongly-typed process interfaces.
643 The rich set of options and features of these tools make them more difficult to learn, but at the
644 same time expose a large number of possibilities.

645 7 Outlook

646 This overview is not meant to be static, but we plan to continue the documentation online in
647 the git repository [16] that contains the implementation of the simple use case. This allows
648 us to take into consideration other WfMSs in the future, and to extend the documentation
649 accordingly. In particular, we would like to make the repository a community effort allowing
650 others to contribute either by modifications of the existing tools or adding new WfMSs. All of
651 our workflow implementations are continuously and automatically tested using GitHub Actions
652 [https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements/ac](https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements/actions)
653 [tions](https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements/actions), which may act as an additional source of documentation on how to launch the workflows.

654 One of the challenges we have identified is the use of container technology in the HPC envi-
655 ronment. In most cases, the way users should interact with such a system is through a module
656 system provided by the system administrators. The module system allows to control the software
657 environment (versioning, compilers) in a precise manner, but the user is limited to the provided
658 software stack. For specific applications, self-written code can be compiled using the available
659 development environment and subsequently run on the system, which is currently the state of
660 the art in using HPC systems. However, this breaks the portability of the workflow.

661 Container technology, employing the “build once and run anywhere” concept, seems to be a
662 promising solution to this problem. Ideally, one would like to be able to run the container
663 application on the HPC system, just as any other MPI-distributed application. Unfortunately,
664 there are a number of problems entailed with this approach.

665 When building the container, great care must be taken with regard to the MPI configuration,
666 such that it can be run successfully across several nodes. Another issue is the configuration of
667 Infiniband drivers. The container has to be build according to the specifics of the HPC system
668 that is targeted for execution. From the perspective of the user, this entails a large difficulty,
669 and we think that further work needs to be done to find solutions which enable non-experts in
670 container technology to execute containerized applications successfully in an HPC environment.

671 Furthermore, challenges related to the joint development of workflows became apparent. In this
672 regard, strongly-typed interfaces are required in order to minimize errors and transparently and
673 clearly communicate the metadata (inputs, outputs) associated with a process in the workflow.
674 This is recommended both for single parameters, but it would be also great to extend that idea
675 to files - not only defining the file type which is already possible within *CWL* - but potentially
676 allowing a type checking of the complete data structure within the file. However, based on our
677 experience with the selected tools, these interfaces and their benefits come at the cost of some
678 form of plugin or wrapper around the software that is to be executed, thus possibly limiting the
679 functionality of the wrapped tool. This means there is a trade-off between easy authoring of the
680 workflow definition (e. g. easily executing any shell command) and implementation overhead
681 for the sake of well-defined interfaces.

682 Another aspect is how the workflow logic can be communicated efficiently. Although each of
683 the tools allows to generate a graph of the workflow, the dependencies between processes can
684 only be visualized for an executable implementation of the workflow, which most likely does
685 not exist in early stages of the project where it is needed the most.

686 An important aspect is the documentation of the workflow results and how they have been
687 obtained. Most tools offer an option to export the data provenance graph, however it would be
688 great to define a general standard supported by all tools as e.g. provided by *CWLProv* [26].

689 **Financial disclosure**

690 None reported.

691 **Conflict of interest**

692 The authors declare no potential conflict of interests.

693 8 Acknowledgements

694 The authors would like to thank the Federal Government and the Heads of Government of the
 695 Länder, as well as the Joint Science Conference (GWK), for their funding and support within the
 696 framework of the NFDI4Ing consortium. Funded by the German Research Foundation (DFG) -
 697 project number 442146713. Moreover, we would like to thank Sebastiaan P. Huber, Michael
 698 R. Crusoe, Eduardo Schettino, Ricardo Wurmus, Paolo Di Tommaso and Johannes Köster for
 699 their valuable remarks and comments on an earlier version of this article and the workflow
 700 implementations.

701 9 Roles and contributions

702 **Philipp Diercks:** Investigation; methodology; software; writing - original draft; writing - review
 703 and editing.

704 **Dennis Gläser:** Investigation; methodology; software; writing - original draft; writing - review
 705 and editing.

706 **Ontje Lünsdorf:** Investigation (supporting); software; writing - original draft (supporting).

707 **Michael Selzer:** Writing - review and editing (supporting).

708 **Bernd Flemisch:** Conceptualization (supporting); Funding acquisition; Project administration;
 709 Writing - review and editing.

710 **Jörg F. Unger:** Conceptualization (lead); Funding acquisition; Project administration; Writing -
 711 original draft (supporting); Writing - review and editing.

712 References

- 713 [1] Enis Afgan et al. “The Galaxy platform for accessible, reproducible and collaborative
 714 biomedical analyses: 2018 update”. In: *Nucleic Acids Research* 46.W1 (May 2018),
 715 W537–W544. ISSN: 0305-1048. DOI: [10.1093/nar/gky379](https://doi.org/10.1093/nar/gky379). eprint: [https://aca-](https://academic.oup.com/nar/article-pdf/46/W1/W537/25110642/gky379.pdf)
 716 [demic.oup.com/nar/article-pdf/46/W1/W537/25110642/gky379.pdf](https://academic.oup.com/nar/article-pdf/46/W1/W537/25110642/gky379.pdf). URL:
 717 <https://doi.org/10.1093/nar/gky379>.
- 718 [2] James Ahrens, Berk Geveci, and Charles Law. “ParaView: An End-User Tool for Large-
 719 Data Visualization”. In: *The Visualization Handbook*. Elsevier, 2005.
- 720 [3] Mohammad Akhlaghi et al. “Toward Long-Term and Archivable Reproducibility”. In:
 721 *Computing in Science & Engineering* 23.3 (May 2021), pp. 82–91. ISSN: 1521-9615,
 722 1558-366X. DOI: [10.1109/mcse.2021.3072860](https://doi.org/10.1109/mcse.2021.3072860). URL: [https://doi.org/10.1109](https://doi.org/10.1109/mcse.2021.3072860)
 723 [/mcse.2021.3072860](https://doi.org/10.1109/mcse.2021.3072860).
- 724 [4] M.S. Alnaes et al. “The FEniCS Project Version 1.5”. In: *Archive of Numerical Software* 3
 725 (2015). DOI: [10.11588/ans.2015.100.20553](https://doi.org/10.11588/ans.2015.100.20553).
- 726 [5] Peter Amstutz et al. *Common Workflow Language, v1.0*. [https://doi.org/10.6084](https://doi.org/10.6084/m9.figshare.3115156.v2)
 727 [/m9.figshare.3115156.v2](https://doi.org/10.6084/m9.figshare.3115156.v2). July 2016. DOI: [10.6084/m9.figshare.3115156.v2](https://doi.org/10.6084/m9.figshare.3115156.v2).

- 728 [6] Arne Babenhauserheide. *SRFI 119: wisp: simpler indentation-sensitive scheme*. <https://srfi.schemers.org/srfi-119/>. June 2015.
- 729
- 730 [7] Oren Ben-Kiki, Clark Evans, and Ingy döt Net. *YAML Ain't Markup Language (YAML) version 1.2*. Accessed: 2022-08-31. Version 1.2. <https://yaml.org/spec/1.2.2/>.
- 731
- 732 2021.
- 733 [8] Michael R. Berthold et al. "KNIME: The Konstanz Information Miner". In: *Data Analysis , Machine Learning and Applications : Proceedings of the 31st Annual Conference of the*
- 734 *Gesellschaft für Klassifikation e. V., Albert-Ludwigs-Universität Freiburg, March 7-9 ,*
- 735 *2007*. New York: Springer, 2007.
- 736
- 737 [9] Neil P. Chue Hong et al. *FAIR Principles for Research Software (FAIR4RS Principles)*.
- 738 <https://doi.org/10.15497/RDA00068>. Version 1.0. May 2022. DOI: [10.15497](https://doi.org/10.15497/RDA00068)
- 739 [/RDA00068](https://doi.org/10.15497/RDA00068). URL: <https://doi.org/10.15497/RDA00068>.
- 740 [10] Iacopo Colonnelli et al. "StreamFlow: cross-breeding cloud with HPC". In: *IEEE Trans-*
- 741 *actions on Emerging Topics in Computing* 9.4 (2021), pp. 1723–1737. DOI: [10.1109](https://doi.org/10.1109/TETC.2020.3019202)
- 742 [/TETC.2020.3019202](https://doi.org/10.1109/TETC.2020.3019202).
- 743 [11] Ludovic Courtès. "Building a Secure Software Supply Chain with GNU Guix". In: *The*
- 744 *Art, Science, and Engineering of Programming* 7.1 (June 2022). ISSN: 2473-7321. DOI:
- 745 [10.22152/programming-journal.org/2023/7/1](https://doi.org/10.22152/programming-journal.org/2023/7/1). URL: [https://doi.org/10.2](https://doi.org/10.22152/programming-journal.org/2023/7/1)
- 746 [2152/programming-journal.org/2023/7/1](https://doi.org/10.22152/programming-journal.org/2023/7/1).
- 747 [12] Ludovic Courtès. "Functional Package Management with Guix". In: *European Lisp*
- 748 *Symposium* (June 2013). DOI: [10.48550/ARXIV.1305.4584](https://doi.org/10.48550/ARXIV.1305.4584). URL: [https://arxiv](https://arxiv.org/abs/1305.4584)
- 749 [.org/abs/1305.4584](https://arxiv.org/abs/1305.4584).
- 750 [13] Michael R. Crusoe et al. "Methods included. standardizing computational reuse and
- 751 portability with the Common Workflow Language". In: *Commun. ACM* 65.6 (June 2022),
- 752 pp. 54–63. ISSN: 0001-0782, 1557-7317. DOI: [10.1145/3486897](https://doi.org/10.1145/3486897). URL: [https://do](https://doi.org/10.1145/3486897)
- 753 [i.org/10.1145/3486897](https://doi.org/10.1145/3486897).
- 754 [14] Ewa Deelman et al. "Pegasus, a workflow management system for science automation".
- 755 In: *Future Generation Computer Systems* 46 (2015), pp. 17–35. ISSN: 0167-739X. DOI:
- 756 [10.1016/j.future.2014.10.008](https://doi.org/10.1016/j.future.2014.10.008). URL: [https://www.sciencedirect.com/sci](https://www.sciencedirect.com/science/article/pii/S0167739X14002015)
- 757 [ence/article/pii/S0167739X14002015](https://www.sciencedirect.com/science/article/pii/S0167739X14002015).
- 758 [15] Paolo Di Tommaso et al. "Nextflow enables reproducible computational workflows".
- 759 In: *Nat Biotechnol* 35.4 (Apr. 2017), pp. 316–319. ISSN: 1087-0156, 1546-1696. DOI:
- 760 [10.1038/nbt.3820](https://doi.org/10.1038/nbt.3820). URL: <https://doi.org/10.1038/nbt.3820>.
- 761 [16] Philipp Diercks et al. *NFDI4Ing Scientific Workflow Requirements*. Version 0.0.1. <https://github.com/BAMresearch/NFDI4IngScientificWorkflowRequirements>. July
- 762 2022.
- 763
- 764 [17] Directorate-General for Research and Innovation (European Commission). *First report*
- 765 *and recommendations of the Commission high level expert group on the European Open*
- 766 *Science Cloud*. Available at <https://op.europa.eu/s/wGAL>. 2016. DOI: [10.2777/9](https://doi.org/10.2777/940154)
- 767 [40154](https://doi.org/10.2777/940154).

- 768 [18] Philip Ewels et al. “Cluster Flow: A user-friendly bioinformatics workflow tool [version 2;
769 referees: 3 approved].” In: *F1000Research* 5 (2016), p. 2824. DOI: [10.12688/f1000res](https://doi.org/10.12688/f1000research.10335.2)
770 [earch.10335.2](https://doi.org/10.12688/f1000research.10335.2). URL: <http://dx.doi.org/10.12688/f1000research.10335.2>.
- 771 [19] Christophe Geuzaine and Jean-François Remacle. “Gmsh: A 3-D finite element mesh
772 generator with built-in pre- and post-processing facilities. THE GMSH PAPER”. In:
773 *Int. J. Numer. Meth. Engng.* 79.11 (May 2009), pp. 1309–1331. ISSN: 0029-5981. DOI:
774 [10.1002/nme.2579](https://doi.org/10.1002/nme.2579). eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.10](https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.2579)
775 [02/nme.2579](https://doi.org/10.1002/nme.2579). URL: <https://doi.org/10.1002/nme.2579>.
- 776 [20] Carole Goble et al. “FAIR Computational Workflows”. In: *Data Intelligence* 2.1-2 (Jan.
777 2020), pp. 108–121. ISSN: 2641-435X. DOI: [10.1162/dint_a_00033](https://doi.org/10.1162/dint_a_00033). eprint: [http](http://direct.mit.edu/dint/article-pdf/2/1-2/108/1893377/dint_a_00033.pdf)
778 [s://direct.mit.edu/dint/article-pdf/2/1-2/108/1893377/dint_a_0003](http://direct.mit.edu/dint/article-pdf/2/1-2/108/1893377/dint_a_00033.pdf)
779 [3.pdf](http://direct.mit.edu/dint/article-pdf/2/1-2/108/1893377/dint_a_00033.pdf). URL: https://doi.org/10.1162/dint_a_00033.
- 780 [21] Lars Griem et al. “KadiStudio: FAIR Modelling of Scientific Research Processes”. In:
781 *Data Science Journal* 21.1 (2022).
- 782 [22] Sebastiaan P. Huber. *aiida-shell*. Version 0.2.0. [https://github.com/sphuber/aiid](https://github.com/sphuber/aiida-a-shell)
783 [a-shell](https://github.com/sphuber/aiida-a-shell). June 2022.
- 784 [23] Sebastiaan P. Huber et al. “AiiDA 1.0, a scalable computational infrastructure for automated
785 reproducible workflows and data provenance”. In: *Sci Data* 7.1 (Sept. 2020). ISSN: 2052-
786 4463. DOI: [10.1038/s41597-020-00638-4](https://doi.org/10.1038/s41597-020-00638-4). URL: [https://doi.org/10.1038/s4](https://doi.org/10.1038/s41597-020-00638-4)
787 [1597-020-00638-4](https://doi.org/10.1038/s41597-020-00638-4).
- 788 [24] Anubhav Jain et al. “FireWorks: A dynamic workflow system designed for high-throughput
789 applications”. In: *Concurrency Computat.: Pract. Exper.* 27.17 (May 2015), pp. 5037–
790 5059. ISSN: 1532-0626, 1532-0634. DOI: [10.1002/cpe.3505](https://doi.org/10.1002/cpe.3505). URL: [https://doi.o](https://doi.org/10.1002/cpe.3505)
791 [rg/10.1002/cpe.3505](https://doi.org/10.1002/cpe.3505).
- 792 [25] Ivo Jimenez et al. “The Popper Convention: Making Reproducible Systems Evaluation
793 Practical”. In: 2017 IEEE International Parallel and Distributed Processing Symposium
794 Workshops (IPDPSW). IEEE, May 2017, pp. 1561–1570. DOI: [10.1109/ipdpsw.2017](https://doi.org/10.1109/ipdpsw.2017.157)
795 [.157](https://doi.org/10.1109/ipdpsw.2017.157). URL: <https://doi.org/10.1109/ipdpsw.2017.157>.
- 796 [26] Farah Zaib Khan et al. “Sharing interoperable workflow provenance: A review of best
797 practices and their practical application in CWLProv”. In: *GigaScience* 8.11 (Nov. 2019).
798 ISSN: 2047-217X. DOI: [10.1093/gigascience/giz095](https://doi.org/10.1093/gigascience/giz095). URL: [https://doi.org/1](https://doi.org/10.1093/gigascience/giz095)
799 [0.1093/gigascience/giz095](https://doi.org/10.1093/gigascience/giz095).
- 800 [27] Johannes Köster and Sven Rahmann. “Snakemake—a scalable bioinformatics workflow
801 engine”. In: *Method. Biochem. Anal.* 34.20 (May 2018), pp. 3600–3600. ISSN: 1367-4803,
802 1460-2059. DOI: [10.1093/bioinformatics/bty350](https://doi.org/10.1093/bioinformatics/bty350). URL: [https://doi.org/10](https://doi.org/10.1093/bioinformatics/bty350)
803 [.1093/bioinformatics/bty350](https://doi.org/10.1093/bioinformatics/bty350).
- 804 [28] Samuel Lampa et al. “SciPipe: A workflow library for agile development of complex
805 and dynamic bioinformatics pipelines”. In: *GigaScience* 8.5 (Apr. 2019). ISSN: 2047-
806 217X. DOI: [10.1093/gigascience/giz044](https://doi.org/10.1093/gigascience/giz044). eprint: [https://academic.oup.c](https://academic.oup.com/gigascience/article-pdf/8/5/giz044/28538382/giz044.pdf)
807 [om/gigascience/article-pdf/8/5/giz044/28538382/giz044.pdf](https://academic.oup.com/gigascience/article-pdf/8/5/giz044/28538382/giz044.pdf). URL:
808 <https://doi.org/10.1093/gigascience/giz044>.

- 809 [29] Soohyun Lee et al. “Tibanna: Software for scalable execution of portable pipelines on the
810 cloud”. In: *Method. Biochem. Anal.* 35.21 (May 2019), pp. 4424–4426. ISSN: 1367-4803,
811 1460-2059. DOI: [10.1093/bioinformatics/btz379](https://doi.org/10.1093/bioinformatics/btz379). eprint: [https://academic
812 .oup.com/bioinformatics/article-pdf/35/21/4424/31617561/btz379.pdf](https://academic.oup.com/bioinformatics/article-pdf/35/21/4424/31617561/btz379.pdf).
813 URL: <https://doi.org/10.1093/bioinformatics/btz379>.
- 814 [30] Felix Mölder et al. “Sustainable data analysis with Snakemake”. In: *F1000Res* 10 (Apr.
815 2021), p. 33. ISSN: 2046-1402. DOI: [10.12688/f1000research.29032.2](https://doi.org/10.12688/f1000research.29032.2). URL:
816 <https://doi.org/10.12688/f1000research.29032.2>.
- 817 [31] Barend Mons et al. “The FAIR Principles: First Generation Implementation Choices and
818 Challenges”. In: *Data Intelligence* 2.1-2 (Jan. 2020), pp. 1–9. ISSN: 2641-435X. DOI:
819 [10.1162/dint_e_00023](https://doi.org/10.1162/dint_e_00023). URL: https://doi.org/10.1162/dint_e_00023.
820 [023](https://doi.org/10.1162/dint_e_00023).
- 821 [32] Simon P. Sadedin, Bernard Pope, and Alicia Oshlack. “Bpipe: A tool for running and man-
822 aging bioinformatics pipelines”. In: *Method. Biochem. Anal.* 28.11 (Apr. 2012), pp. 1525–
823 1526. ISSN: 1460-2059, 1367-4803. DOI: [10.1093/bioinformatics/bts167](https://doi.org/10.1093/bioinformatics/bts167). eprint:
824 [https://academic.oup.com/bioinformatics/article-pdf/28/11/1525/1690
825 5290/bts167.pdf](https://academic.oup.com/bioinformatics/article-pdf/28/11/1525/16905290/bts167.pdf). URL: <https://doi.org/10.1093/bioinformatics/bts167>.
- 826 [33] Michael A. Salim et al. *Balsam: Automated Scheduling and Execution of Dynamic, Data-
827 Intensive HPC Workflows*. <https://arxiv.org/abs/1909.08704>. 2019. DOI:
828 [10.48550/ARXIV.1909.08704](https://arxiv.org/abs/1909.08704). URL: <https://arxiv.org/abs/1909.08704>.
- 829 [34] Eduardo Naufel Schettino. *pydoit/doi: Task management & automation tool (python)*.
830 <https://doi.org/10.5281/zenodo.4892136>. June 2021. DOI: [10.5281/zenodo
831 .4892136](https://doi.org/10.5281/zenodo.4892136). URL: <https://doi.org/10.5281/zenodo.4892136>.
- 832 [35] Nico Schlömer. *meshio: Tools for mesh files*. [https://doi.org/10.5281/zeno
833 do.6346837](https://doi.org/10.5281/zenodo.6346837). Version v5.3.4. Mar. 2022. DOI: [10.5281/zenodo.6346837](https://doi.org/10.5281/zenodo.6346837). URL:
834 <https://doi.org/10.5281/zenodo.6346837>.
- 835 [36] Will Schroeder et al. *The visualization toolkit : an object-oriented approach to 3D graphics*.
836 4th ed. Kitware, 2006.
- 837 [37] Michael Sperber et al. “Revised6 Report on the Algorithmic Language Scheme”. In: *J.
838 Funct. Program.* 19.S1 (Aug. 2009), p. 1. ISSN: 0956-7968, 1469-7653. DOI: [10.1017
839 /s0956796809990074](https://doi.org/10.1017/s0956796809990074). URL: <https://doi.org/10.1017/s0956796809990074>.
- 840 [38] Martin Uhrin et al. “Workflows in AiiDA: Engineering a high-throughput, event-based
841 engine for robust and modular computational workflows”. In: *Nato. Sc. S. Ss. Iii. C. S.* 187
842 (Feb. 2021), p. 110086. ISSN: 0927-0256. DOI: [10.1016/j.commatsci.2020.110086](https://doi.org/10.1016/j.commatsci.2020.110086).
843 URL: <https://doi.org/10.1016/j.commatsci.2020.110086>.
- 844 [39] John Vivian et al. “Toil enables reproducible, open source, big biomedical data analyses”.
845 In: *Nat Biotechnol* 35.4 (Apr. 2017), pp. 314–316. ISSN: 1087-0156, 1546-1696. DOI:
846 [10.1038/nbt.3772](https://doi.org/10.1038/nbt.3772). URL: <https://doi.org/10.1038/nbt.3772>.
- 847 [40] Kate Voss, Geraldine Van Der Auwera, and Jeff Gentry. *Full-stack genomics pipelining
848 with GATK4 + WDL + Cromwell [version 1; not peer reviewed]*. slides. [https://f1000
849 research.com/slides/6-1381](https://f1000research.com/slides/6-1381). 2017. DOI: [10.7490/f1000research.1114634.1](https://doi.org/10.7490/f1000research.1114634.1).
850 URL: <https://f1000research.com/slides/6-1381>.

- 851 [41] Mark D. Wilkinson et al. “The FAIR Guiding Principles for scientific data management
852 and stewardship”. In: *Sci Data* 3.1 (Mar. 2016). ISSN: 2052-4463. DOI: [10.1038/sdat](https://doi.org/10.1038/sdata.2016.18)
853 [a.2016.18](https://doi.org/10.1038/sdata.2016.18). URL: <https://doi.org/10.1038/sdata.2016.18>.
- 854 [42] Peter Williams and Contributors. *The Tectonic Typesetting System*. [https://tectonic-](https://tectonic-typesetting.github.io/en-US/)
855 [typesetting.github.io/en-US/](https://tectonic-typesetting.github.io/en-US/). Accessed: 2022-06-02. 2022.
- 856 [43] Andy Wingo et al. *GNU Guile*. <https://www.gnu.org/software/guile/>. Feb. 2022.
- 857 [44] Ricardo Wurmus et al. *GUIX Workflow Language*. <https://guixwl.org>. Version 0.5.0.
858 July 2022.
- 859 [45] Andy B. Yoo, Morris A. Jette, and Mark Grondona. “SLURM: Simple Linux Utility for
860 Resource Management”. In: *Job Scheduling Strategies for Parallel Processing*. Ed. by
861 Dror Feitelson, Larry Rudolph, and Uwe Schwiegelshohn. Berlin, Heidelberg: Springer
862 Berlin Heidelberg, 2003, pp. 44–60. ISBN: 978-3-540-39727-4.