

1 **Model coupling through reproducible adapter workflows based on shared transformation** 2 **functions**

3 Patrick Kuckertz^{1,4,*}, Benjamin Fuchs², Julian Schönau¹, Hedda Gardian², Kevin Knosala¹,
4 Eugenio Arellano Ruiz², Jan Göpfert^{1,3}, Hans Christian Gils², Jann M. Weinand¹, Patrick
5 Jochem², Jochen Linßen¹, and Detlef Stolten^{1,3}

6 ¹Forschungszentrum Jülich GmbH, Institute of Energy and Climate Research – Jülich Systems
7 Analysis, 52425 Jülich, Germany

8 ²German Aerospace Center (DLR), Institute of Networked Energy Systems, 70563 Stuttgart,
9 Germany

10 ³RWTH Aachen University, Chair for Fuel Cells, Faculty of Mechanical Engineering, 52062
11 Aachen, Germany

12 ⁴Lead contact

13 *Correspondence: p.kuckertz@fz-juelich.de

14 **Summary**

15 In computational science, coupling research software and data into workflows is essential for
16 addressing complex research questions and ensuring reproducibility. However, current meta-
17 data schemas rarely provide sufficient information about data models in software interfaces and
18 datasets, hindering their effective integration into workflows. Additionally, most workflow man-
19 agement tools are not designed to handle metadata for process reproduction. This article refines
20 the DataDesc metadata schema to annotate data models not only in software interfaces but also
21 in datasets, and presents an extension to the DataDesc framework for automatically compar-
22 ing data models and identifying transformation requirements. The ioProc workflow manager is
23 introduced to bundle shared transformation functions into adapter workflows and ensure trans-
24 parent process documentation. Two use cases from energy systems analysis demonstrate the
25 workflow design and implementation approaches. Overall, this article bridges the gap between
26 abstract guidelines, such as the FAIR principles, and researchers' daily data and software driven
27 analyses, promoting reusability, reproducibility and transparency in science.

28 **Keywords**

29 Software Coupling, Data Integration, Scientific Workflows, Data Models, Transformation Func-
30 tions, Metadata, Reusability, Interoperability, DataDesc Schema, ioProc Workflow Manager

31 **Introduction**

32 **Motivation**

33 In computational sciences and especially in energy systems analysis, the coupling of research
34 software into comprehensive software workflows is a common approach for incorporating multi-
35 ple perspectives into a scientific analysis and to investigate complex research questions.¹ In a
36 workflow context, the sustainable handling of resources is particularly important in two aspects.
37 On the one hand, these workflows represent in themselves flexible research tools. Reuse or
38 adaption of these workflows to changes in research questions in a fast, reliable and low-effort
39 approach can avoid considerable redundant development, implementation, and maintenance

40 efforts but requires careful upfront design and additional effort.² On the other hand, these work-
41 flows play a central role in scientific exchange, as they contribute to the transparency and relia-
42 bility of scientific conclusions by enabling others to trace, interpret and reproduce the results of
43 these workflows.^{3,4} Workflow design therefore has the potential to positively increase the overall
44 efficiency of a research community.

45 Usually, however, scientists are confronted with major challenges when creating and applying
46 workflows that should be adaptable, reuseable and reproducible by others. Especially missing or
47 difficult-to-process information about software and data resources represents a major obstacle.
48 Particularly in a research field such as energy systems analysis, in which no uniform interface
49 and data standards have been established and which is characterized by heterogeneous for-
50 mats originating from different research disciplines, the provisioning of metadata which includes
51 information about individually used data models is imperative.⁵ Another hindrance is the lack of
52 computer-aided workflow design and building tools that cover and support both research soft-
53 ware and data. In an environment full of initially incompatible data models, there is a strong need
54 for reliable and robust data transformation and processing capabilities, in order to combine input
55 data and software models into a structured process. As long as there are no adequate tools sup-
56 porting scientists in workflow integration or offering automation capabilities, data models must be
57 manually researched and compared, which is a tedious and time-consuming task. In addition,
58 transformation requirements must be identified manually and suitable transformation functions
59 researched in a time-consuming process. A third challenge lies in the lack of traceability and
60 reproducibility of workflows and workflow results. Even as there is no lack of workflow tools with
61 which researchers can implement their model chains (see for example Kieser et al.⁶ or Mölder
62 et al.⁷), the metadata accompanying these processes, which is relevant for their understanding
63 and interpretation, must be recorded, linked and stored separately, usually in a manual process.
64 This additional effort, which has to be made for each individual model coupling, data processing
65 or transformation, is often not done consistently. This situation contributes to the current repro-
66 ducibility crisis in science.⁸⁻¹⁰ An illustrative example is a study of Jupyter Notebooks, a common
67 form of documenting scientific workflows and analysis, for their reproducibility. The study of Pi-
68 mentel et al.¹¹ showed issues with dependency documentation and execution of the notebooks
69 and thus a lack of reproducibility, transparency and adaptability.

70 Confronted with this situation, a cultural change is slowly taking place in energy systems
71 analysis and other science communities, which places increasing importance on overcoming the
72 associated challenges. In particular, the objectives formulated in the FAIR principles, findability,
73 accessibility, interoperability and reusability, are effectively promoting the adoption of metadata
74 standards as this directly contributes to all four principles.¹² At the same time, the joint devel-
75 opment and use of domain ontologies for the uniform description of software, data, and other
76 research data artifacts, such as the Open Energy Ontology¹³, is also gaining traction. As the
77 progress of adopting the FAIR principles continues, numerous problems and uncertainties in
78 the practical implementation of these guidelines are experienced in everyday science. To solve
79 these, novel approaches and solutions are required with renewed urgency. As described by
80 Leipzig et al.¹⁴ there exists an explicit need for action in the area of metadata description of data
81 that is to be used as input by software models. The authors call for a common metadata stan-
82 dard that enables the description of file formats and data structures at workflow level and their
83 automated comparison in order to link individual elements of analytical processes. They claim
84 that metadata tools are as important for practical and computational research as the software
85 and data itself. Leipzig et al.¹⁴ further show that despite the existence of numerous workflow
86 tools, collecting and analyzing provenance, i.e. recording all activities that lead to the creation of
87 a data object, is still a key challenge for workflow design.

88 The approaches presented in this article provide an answer to the question how metadata
89 aware workflows, which cover software and data, can be implemented and documented to make

90 a significant step towards the FAIR principles. In the *Related Work* section, current metadata
91 schemas are first compared and examined with regard to their suitability for describing data
92 models. Subsequently, an introduction to data comparison and workflow tools is given, focusing
93 in particular on their capabilities for metadata processing in a scientific context. The *Results*
94 section shows the extent to which the metadata elements of the DataDesc schema, which is
95 designed to describe data models in software interfaces, are also suitable for describing data
96 models in data files, and which structural refinements have been applied to facilitate this applica-
97 tion case. Furthermore, an extension of the DataDesc framework for the automated comparison
98 of metadata is presented, in which data models are compared and any discrepancies are identi-
99 fied as transformation needs, whose machine-actionable description can be used in the future
100 for the automated identification and reuse of data transformation methods in modular workflow
101 designs. In a second step, the ioProc workflow manager¹⁵ is presented, with which data trans-
102 formation steps can be combined into so-called adapters as needed and sustainably reused in
103 workflows. The individual process steps are also automatically documented as metadata. Fi-
104 nally, the *Results* section presents two use cases from the context of energy systems analysis
105 to demonstrate how DataDesc and ioProc can be used to design and implement transparent,
106 reproducible and well-documented scientific workflows. The *Discussion* section, which empha-
107 sizes the main features of the approaches presented as well as their limitations, provides an
108 insight into future work and closes this article.

109 This work hence presents an adaptable and extensible approach for an improved design and
110 reusability of scientific workflows incorporating software and data, rooted in and exploiting the
111 benefits of metadata. The presented approach contributes to the FAIR principles and supports
112 the transparency and reproducibility of workflows in a scientific community by adding automatic
113 metadata generation.

114 **Related Work**

115 This section examines current metadata schemas in regards to their suitability for describing
116 data models. Furthermore, different approaches for comparing data are contrasted. Finally, the
117 role of workflow tools for computational sciences is established.

118 **Metadata Schemas**

119 A metadata schema standardizes the description of artifacts within its scope by defining a set
120 of metadata elements to be used.¹⁶ A metadata standard may further standardize the encoding
121 format, allowed values, their representation, and so forth.¹⁶ Application profiles adapt subsets
122 of one or multiple metadata schemas to tailor these to specific applications. They may further
123 define custom elements, rules, best practices, etc. Many schemas, standards and application
124 profiles exist with different scopes. The scope refers to specific use cases that are supported,
125 but may be limited to, for example, a particular artifact type (such as images, software, datasets,
126 etc.), scientific domain, or country.

127 Domain-agnostic standards that propose sets of elements to describe datasets include Dublin
128 Core¹⁷, DataCite¹⁸ that provides a mapping to Dublin Core, DCAT¹⁹ that incorporates terms from
129 Dublin Core as well as other controlled vocabularies, and schema.org's Dataset class²⁰ that in
130 turn is based on DCAT. None of these describe data models in detail. In addition, the general
131 metadata of a dataset such as *authors*, *titles*, *descriptions*, the content and how to access it
132 is only defined in an abstract way using terms such as *encoding*, *file format*, *has part*, *about*,
133 *variable measured*, *temporal or spatial resolution and coverage*. Dublin Core and DCAT allow to

reference a schema or standard using the property *conforms to*. However, there is no schema describing the data model provided.

RO-Crate²¹ and Data Package²² are approaches to bundling data and their metadata. For example, Data Package allows the specification of table schemas for individual files. OMETADATA²³ is a domain-specific standard that builds on Data Package and encourages the referencing of concepts from controlled vocabularies and ontologies.

Representing n-dimensional datasets in the Resource Description Framework (RDF), the RDF Data Cube Vocabulary²⁴ is not limited to tabular data. It represents the actual data itself, but uses the rich semantics of RDF to describe information that would normally be part of metadata. D-REPR²⁵ and the Software Description Ontology²⁶ reuse the RDF Data Cube Vocabulary to define the structure of a dataset. Other self-describing data formats include HDF5²⁷ and NetCDF²⁸.

More recently, Croissant²⁹ and DataDesc³⁰ proposed approaches that are independent of the dataset format, allowing for the referencing of terms from controlled vocabularies, and the description of multidimensional as well as nested data structures. To increase the reusability of datasets within the machine learning community, Croissant supports the specification of data types and foreign key references. In addition, usage information, e.g. where the training, validation and test split is located, which columns are to be extracted and how the values are to be parsed, facilitates the training of machine learning models. DataDesc is a domain-agnostic standard for describing software interfaces and their data models with metadata to facilitate model coupling, workflow composition, and software discovery as well as integration in general. It builds on the OpenAPI specification, a widely used standard to document APIs, and provides a mapping to schema.org. DataDesc offers the possibility to describe multi-dimensional and nested function parameter structures and value ranges in runtime representations. Semantic concepts of variables as well as their units or quantity types can be described in a machine-actionable way by referring to terms from controlled vocabularies.

In its area of application for interface documentation, the DataDesc schema offers the highest level of detail and the greatest flexibility in the description of data models, so that it can map variables and function parameters of any complexity. In addition, it integrates and references current metadata standards and can be used across domains and applications. However, there is currently no metadata schema that has been designed to map data models of datasets at a comparable level of detail.

Data Comparison Tools

To what extent and in what respect data is compared depends on the use case. We focus on the comparison of data models between datasets and software or service interfaces with the aim of integrating them into comprehensive software workflows. This requires consideration of both semantics and technical representations, including data structures, data formats, and value range specifications.

To compare different versions of files, utilities such as GNU Diffutils³¹ or Git's diff³² highlight differences based on the line-based edit distances between them. With binary files, a comparison of lines or other text chunks is generally not meaningful; they are either classified as different or not.³¹ Specifying data models in separate metadata files circumvents this limitation.

To compare the contents of files beyond the purpose of comparing versions, the similarity of individual strings, such as the keys and values in a metadata file, can be calculated based on their token- or character-based edit distance (e.g., using the Levenshtein distance³³). Often strings are normalized beforehand (e.g., by lowering and lemmatizing of words). To also consider semantics, the vector similarity between strings can be calculated based on word embeddings (such as Word2Vec³⁴ and GloVe³⁵ embeddings). Alternatively, strings can be mapped

182 to concepts in ontologies or lexical databases (such as WordNet³⁶) and compared based on
183 the semantic relations they define, such as a synonyms or hyponyms. The meaning of a word
184 or phrase depends on its context. Hence, consideration of the context is important. A-Match³⁷
185 is an example of an application that uses both string metrics and ontologies to compare APIs.
186 Following a human-in-the-loop approach, the system's suggestions are visualized in a GUI and
187 can be accepted or rejected.

188 When comparing more complex entities than strings, more information can be considered.
189 For example, in addition to the parameter name, the allowed or actual value range can be con-
190 sidered (are the values within the expected range?) as well as the units (can they be converted
191 into each other?). Consideration of the data structures is particularly important when compar-
192 ing dataset and software interface descriptions. For RDF graphs, their structure, value ranges
193 and other constraints can be defined using the Shapes Constraint Language (SHACL)³⁸. Corre-
194 sponding validators can be used to validate RDF graphs based on these definitions. The MINT
195 (Model INTEgration) framework³⁹ also requires the mapping of datasets to RDF. Variables are
196 associated with a Scientific Variables Ontology (SVO) concept. Based on the unified RDF rep-
197 resentations, both content and structure are considered when matching a set of user-selected
198 target variables to models that can output them.

199 When comparing data models between datasets and software or services, both semantics
200 and data structures must be taken into account, ideally independent of the file format to enable
201 broad applicability. The functional extension of the DataDesc framework presented in this work
202 encourages mapping to controlled vocabularies, is not restricted to a particular file format, and
203 does not require the conversion of heterogeneous datasets into a unified representation.

204 **Workflow Tools**

205 Workflow tools play an important role in computational scientific work.⁷ The use of workflow tools
206 ranges from high performance computing environments to individual data analysis pipelines of
207 scientists. Software solutions are correspondingly diverse, with software such as SLURM⁴⁰
208 in HPC contexts not only managing workflows but integrating resource management down to
209 Jupyter notebooks⁴¹ in individual data analysis pipelines with loose but easily adaptable work-
210 flows. The side-benefit for science of such workflow tools, especially when they provide a static
211 workflow definition, is the creation of artefacts that are themselves suitable for documentation,
212 and the further creation of traceability and reproducibility for these workflows. A widely used
213 workflow tool is Snakemake⁴², which binds programs to flexible, statically declared workflows
214 and provides a reproducible execution environment. An example of one of the more prominent
215 applications of such tools, at least in energy systems analysis, is the coupling of models into
216 complex workflows for multi-perspective scientific analysis.⁴³

217 Most of the workflow tools used in the scientific community, Snakemake being a rare excep-
218 tion, were not originally developed specifically for scientific activities. As a result, few workflow
219 tools incorporate scientific metadata processing and handling by design. The importance of
220 scientific metadata is now being increasingly recognised.^{44,45} There, the tracking of data-related
221 metadata and the computer-aided generation of metadata for generated datasets is central to the
222 reusability and traceability of scientific results. The increasing popularity of the FAIR principles
223 both among scientists⁴⁴ and funding agencies⁴⁶⁻⁴⁸ only adds to the importance. The support
224 of metadata information during workflows and of dedicated interfaces to scientific infrastructures
225 such as metadata hubs is usually not natively covered by traditional generic workflow tools. Fi-
226 nally, workflow tools introduce their own complexity and require additional skills to be acquired
227 by researchers. As workflow tool skills are a secondary concern in the scientific work environ-
228 ment, workflow managers with easy-to-learn principles and easy-to-use approaches are usually
229 favored.

230 In summary, workflow tools often lack native support for scientific metadata and integration
231 with scientific infrastructures, which makes data traceability and interoperability more difficult.
232 Furthermore, their complexity and steep learning curve make them less attractive in scientific
233 environments where ease of use is a priority.

234 In this work we demonstrate with our implementation of a scientific workflow in ioProc, how
235 the issue of lock-in effects into a specific workflow manager and that workflow tools can be
236 designed, such, that scientists can apply them with their preexisting knowledge in scripting and
237 notebook environments.

238 Results

239 This section first describes a modular approach to software and data integration and shows
240 how model workflows can be practically designed from reusable software components using
241 the refinements and extensions of the DataDesc framework presented here. In a next step, io-
242 Proc is introduced as a workflow management tool that is designed to meet the special needs of
243 scientific research and with which model workflows can be implemented transparently and repro-
244 ducibly. Finally, DataDesc and ioProc are used to demonstrate the design and implementation
245 of two application cases from the domain of energy systems analysis.

246 Designing Reusable Workflows with DataDesc

247 This work takes up the modular workflow concept presented by Kuckertz et al.⁴⁹, which propa-
248 gates a flexible approach to data and software integration and thus offers a lightweight alterna-
249 tive to the laborious introduction and adoption of static interface and data standards in energy
250 systems research (cf. Figure 1). Along the workflow concept, information transfers between soft-
251 ware modules via non-persistent in-memory data formats, such as integers, arrays or classes,
252 and is exchanged between software and data files based on persistent data formats, such as
253 Excel, XML or NetCDF4. This work addresses the integration of persistent data formats into
254 software interfaces. To this end, the concept first requires the metadata description of software
255 and data artifacts, whereby in particular their inherent data models are annotated in detail. Once
256 these are available, the data models mapped by software interfaces can be compared with those
257 used in data files in a second step. This provides information about the compatibility of data and
258 interfaces and identifies transformation requirements that can be covered by data processing.
259 Along this process, the concept aims at the reuse of software and data artifacts and their flexible
260 integration into complex application-specific model workflows. In addition to the metadata de-
261 scriptions, transformation functions form a central component, as they act as a link between the
262 individual software interfaces and heterogeneously structured datasets within the domain of en-
263 ergy research. The concept intends for the transformation functions to be designed in a modular
264 way, described transparently and made available as open-source so that they can be jointly used
265 and sustainably developed within the research community. Overall, the concept is intended to im-
266 prove the reuse of research software and also the reproducibility and validation of study results,
267 thereby promoting efficiency and academic exchange in the field of energy systems analysis.

268 In order to implement this concept in practice, metadata schemas are required that can be
269 used to describe data models in detail both in interfaces and in datasets. These must be com-
270 patible with each other insofar as the annotations created on their basis allow them to be com-
271 pared exactly with each other. In addition, the schemas should support semantic references
272 to ontologies, such as the Open Energy Ontology⁵⁰ developed for energy systems research,
273 and machine-actionable exchange formats so that the comparisons can be carried out automat-

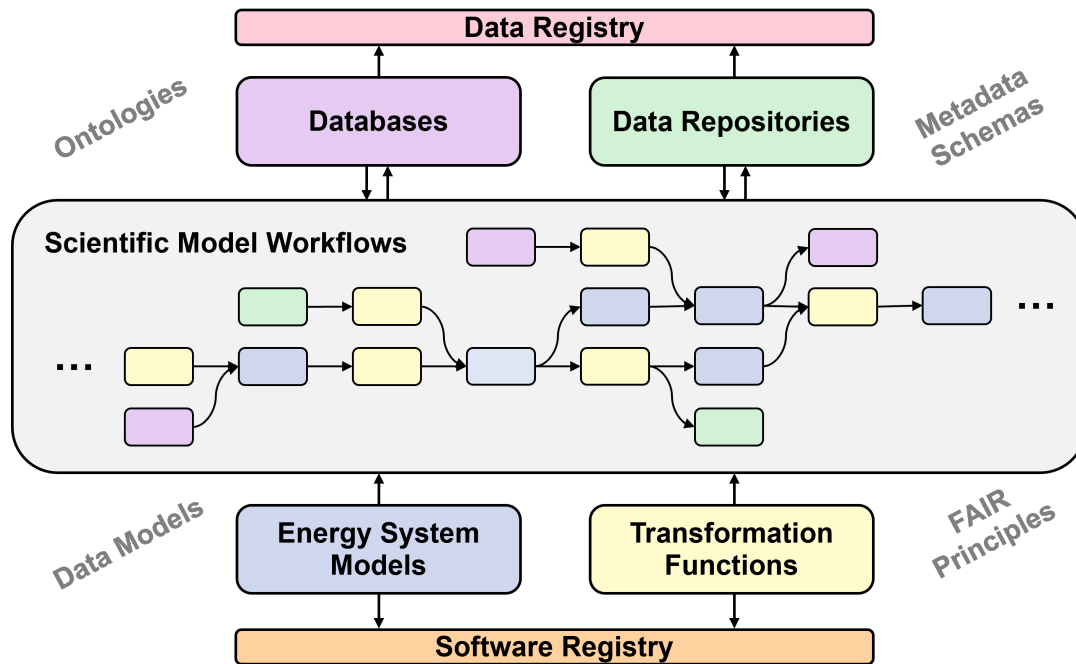


Figure 1: A modular workflow concept for the flexible integration of software and data in the context of a distributed research data infrastructure. The original figure by Kuckertz et al.⁴⁹ was supplemented by a software metadata registration component to facilitate the findability and reuse of energy systems models and data transformation functions.

274 ically. With DataDesc, Kuckertz et al.³⁰ have already developed a framework that meets these
 275 requirements, at least for the formal description of software interfaces. DataDesc centers around
 276 a software metadata schema that describes the data models on which software interfaces are
 277 based. In addition, DataDesc effectively promotes the FAIRness, i.e., findability, accessibility,
 278 interoperability, and reusability, of research software. As described in the *Related Work* section,
 279 there is currently no metadata schema for datasets that can be used to map their data models
 280 at a comparable level of detail. In order to present an approach to compensate for this omission,
 281 it is shown in the following to what extent the elements of the DataDesc metadata schema are
 282 also suitable for describing data models of datasets, and how the schema structure has been
 283 refined for this purpose.

284 The part of the DataDesc schema that is the focus of this work is the *Data Schema Ob-*
 285 *ject* (cf. Figure 2). Its properties can be used to describe even deeply nested data models in
 286 terms of contents, formats, value ranges and structures. The contents of data structures can be
 287 documented using *description* and clearly assigned to individual ontology classes using *seman-*
 288 *ticConcept*. It is of no consequence whether these contents are expected by a software interface
 289 or provided by a data file. The properties *type*, *format*, *mediaType*, and *charSet* can be used to
 290 describe data model components of both persistent and non-persistent variables. For example,
 291 the type and format can be used to specify not only the valid file type for variables that expect a
 292 file as input, but also the assumed data structure within a transferred file. The various *minimum*
 293 and *maximum* properties can be used to define value ranges that must be adhered to when
 294 using an interface in order to ensure error-free data processing. With regard to files, these prop-
 295 erties describe specific value ranges contained in the dataset. Finally, the description of nested,
 296 grouping or dimensionally resolving data structures with *items* and *properties* remains the same,
 297 regardless of whether they are used in a function variable of complex type or a data file. Overall,
 298 it becomes evident that the Data Schema Object can be used to describe data models in files in
 299 the same way as in software interfaces.

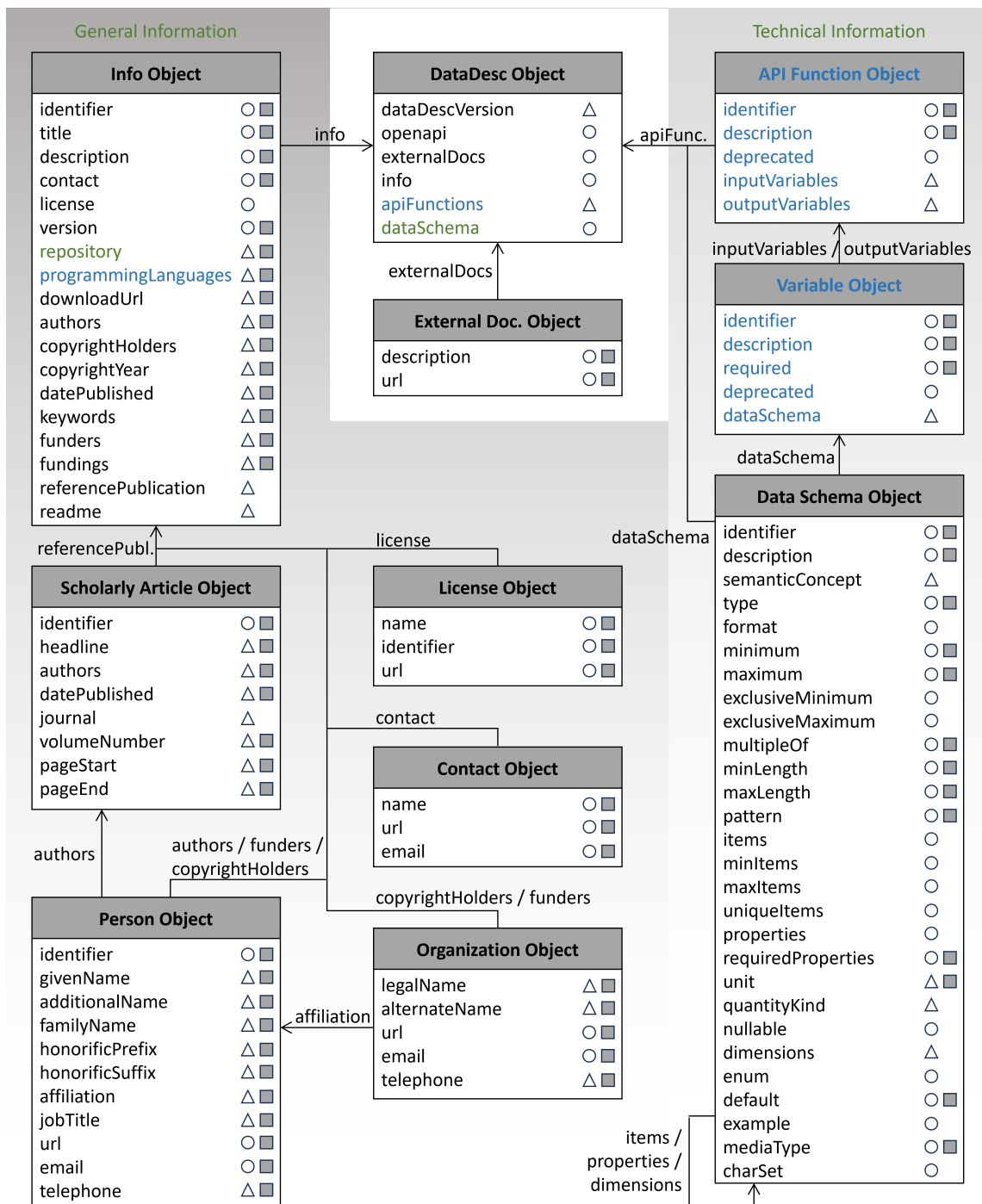


Figure 2: Structure and content of the DataDesc schema for describing software and their interface data models, adapted from Kuckertz et al.³⁰. Both the general and the technical information are organized in information objects, with arrows indicating the various relationships between them. DataDesc properties that map directly or via extensions with the OpenAPI specification are indicated by white circles and white triangles, respectively. Properties that map to the Schema.org ontology are indicated by a gray square. The schema and figure from Kuckertz et al.³⁰ have been slightly adapted by the changes shown in green to accommodate the description of data models mapped in files. Components written in blue are intended exclusively for the description of interfaces, while all other objects and properties are also used to describe datasets. Individual property definitions can be viewed at GitHub: https://github.com/FZJ-IEK3-VSA/DataDesc/blob/main/schema/DataDesc_schema_v1.2.md.⁵¹

300 In Figure 2, object types and properties of the DataDesc schema that are intended exclu-
301 sively for the description of interfaces are highlighted in color. It is clear to see that the majority,
302 including the part required for the description of data models, is just as suitable for the annotation
303 of data files as for software interfaces. The only structural difference is that the *Data Schema Ob-*
304 *ject* is assigned directly to the *DataDesc Object* when describing files. By reusing the schema,
305 also the exchange format of the DataDesc framework can be reused, which inherently ensures
306 the homogeneity of the metadata descriptions.

307 On this basis, the DataDesc framework has been expanded to include analysis functionality
308 that enables the direct and systematic comparison of a data model mapped in a file with the
309 one implemented in a software interface and reveals the degree of their compatibility. In ad-
310 dition to the structures of the respective data models, also the properties of the variables they
311 contain, such as dimensionality, value ranges and units, are contrasted individually. During the
312 comparison process, the software interface is taken as the reference point from which the ex-
313 tent to which the data model of the dataset corresponds to the necessary interface structures is
314 assessed. While the absence of expected variables or deviations from expected units or value
315 ranges reduce the degree of compatibility, additional information contained in the file that is not
316 required by the interface is ignored and does not have a negative effect. Essentially, the compar-
317 ison process checks whether the information required by the software interface is available as a
318 subset in the data file.

319 As a result of a comparison, a report is generated that provides detailed information about
320 (in)compatibilities. Starting from an overall result that reflects whether the comparison has re-
321 vealed discrepancies or not, the report goes into more and more detail along the hierarchy of
322 the data model until it becomes clear at which points in the data model which requirements have
323 not been met. This accurate error indication is intended to support researchers in the design
324 and reuse of computational workflows, to avoid errors in data processing and to make software
325 and data integration more efficient. At the same time, reports follow the DataDesc structure so
326 that they can formally represent compatibility information as independent, machine-actionable
327 DataDesc documents and can be further processed automatically.

328 If a file is to be used as input by a software despite incompatibilities, transformation require-
329 ments arise that must be covered by data processing. Thereby, necessary transformations can
330 vary considerably in their complexity and scope, ranging from simple unit conversions to exten-
331 sive data restructuring. As the report provides information about the type of incompatibility, it
332 can be used to identify suitable transformation functions. For example, differences in currency
333 formats can be resolved with CuCoPy⁵², while discrepancies in the spatial or temporal resolution
334 of the data can be addressed with Spagat⁵³ or tsam^{54,55}.

335 Many such transformation functions are currently available open-source. However, even if
336 they are registered in software catalogs as indicated in Figure 1, they cannot yet be automat-
337 ically suggested or selected on the basis of identified transformation requirements. With their
338 formalized information on data models and incompatibilities, the machine-actionable DataDesc
339 reports form a technical basis for the goal-oriented findability and reusability of data processing
340 software.

341 Once the model and data components required for an analysis have been connected using
342 the transformation functions identified via DataDesc, the modular concept subsequently requires
343 them to be combined into transparent and reproducible computational workflows.

344 **Implementing Reproducible Workflows with ioProc**

345 ioProc is a library-level workflow manager, which distinguishes it from other tools such as the
346 popular Snakemake workflow management system⁷. Snakemake operates on the software level,

347 so it ties individual software pieces together into a larger workflow. ioProc on the other hand
348 chains functions, called actions in this context, into reproducible workflow specified by static
349 definition file, called adapters. This adapter also serves as part of the scientific documentation
350 of the workflow and can be used to run the workflow repeatedly and reproduce analysis results.
351 In addition to the workflow itself, ioProc also writes log files, with one dedicated only to data
352 modifications to adhere to good scientific practices. All read and write operations to and from
353 data structures inside of ioProc actions are thus documented and traceable for each execution
354 of the workflow. Furthermore, ioProc follows the open-closed principle, thus that it is open to
355 extensions but closed to modifications. This manifests in different aspects of the software. For
356 example, all actions are saved in an action folder, including the default set of built-in actions,
357 that ioProc can generate on request. The ioProc generated actions are standardized, but as
358 they are stored as a module, accessible by the user and outside of the ioProc software, opens
359 them up for modification and extension without the need to modify ioProc itself. The user can
360 furthermore declare own actions in additional files. ioProc, if configured with the location of these
361 extension files, can parse these actions and apply them in a user specified workflow. In addition,
362 users can share their actions with others so that they can use them in their own workflows. To
363 improve reuseability, files with ioProc compatible actions can be modified with little mockup code,
364 to make them independent from an ioProc installation. It is then possible to use these actions like
365 ordinary Python functions in environments without ioProc, like a framework, a jupyter notebook
366 or another workflow tool. ioProc thus creates an environment which is geared towards good
367 scientific practices with a unique focus on library level workflows.

368 **Application Cases**

369 To illustrate our design and implementation of reproducible software workflows from reusable
370 software and data components and to show the general applicability of the approaches pre-
371 sented, two independent and simplified model workflows are compared in this section (see Fig-
372 ure 3). With REMix (Renewable Energy Mix)⁵⁶ and ETHOS.FINE (Framework for Integrated En-
373 ergy Assessment (FINE) of the Energy Transformation Pathway Optimization Suite (ETHOS)),^{57,58}
374 two mathematical optimization frameworks with a similar scope from the field of energy systems
375 research are employed. Both application cases are based on pre-existing models that are instan-
376 tiated from the two frameworks. While both models contain individual sets of predefined param-
377 eters and data, they also share some input datasets. The shared input data includes unresolved
378 data regarding energy technologies on the one hand and temporally and spatially resolved re-
379 newable energy potentials on the other. During data processing within the two workflows, these
380 differently structured datasets are adapted to the data models expected by the frameworks' inter-
381 faces along individual transformation steps, whereby some of the modular transformation meth-
382 ods are reused in both workflows. While the REMix workflow is implemented using ioProc and
383 maps a file-based import of the data, the ETHOS.FINE workflow is implemented using Jupyter
384 Notebook on the basis of a script-based data import.

385 When selecting the application cases, the focus was on the fact that they differ in a number
386 of relevant characteristics, but remain structurally comparable. The model calculations carried
387 out are realistic in terms of data structures, formats, contents and value ranges. Beyond that,
388 however, they do not represent conclusive analyses of energy systems.

389 **Shared Input Datasets**

390 Both workflows ingest technology catalogues from the Danish Energy Agency⁵⁹, which provide
391 a comprehensive but unresolved dataset of techno-economic data for technologies commonly

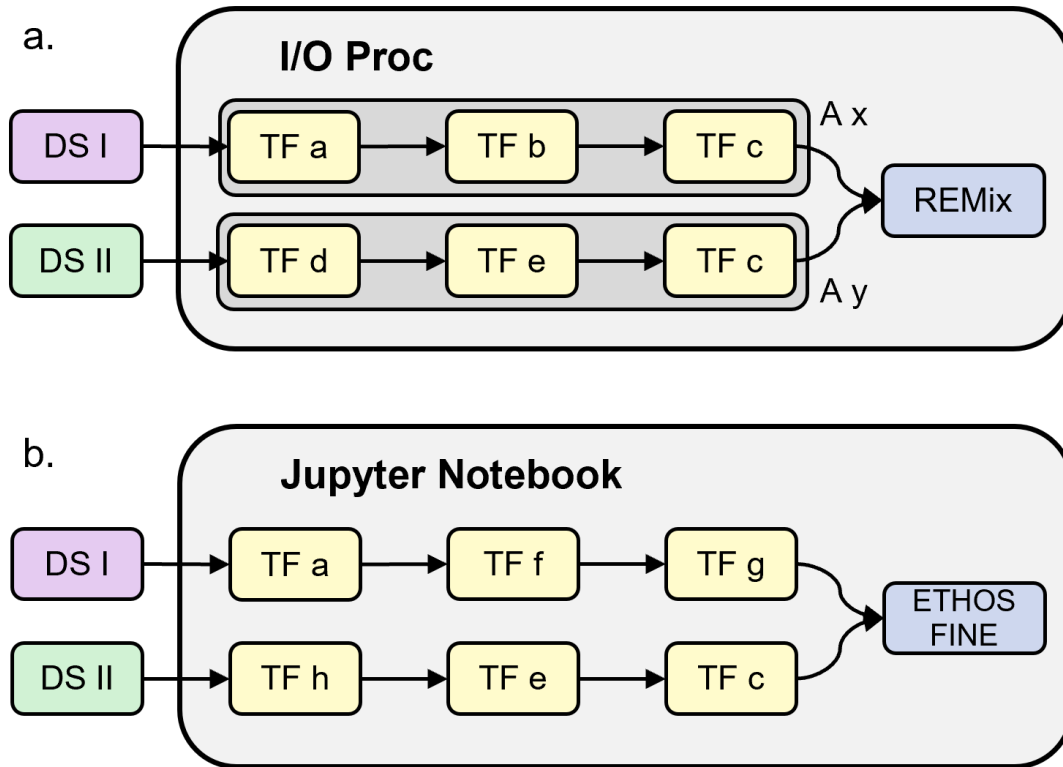


Figure 3: Schematic comparison of an ioProc (a.) and a Jupyter Notebook software workflow (b.) that use the same datasets (DS) as input and transfer their contents into the energy system models REMix and ETHOS.FINE using individually combined transformation functions (TF) and adapters (A).

392 used in energy systems analysis. The data is provided under the CC BY 4.0 license in Excel
 393 spreadsheets grouped by technology type, such as electricity and heat generation⁶⁰, storages⁶¹
 394 and renewable fuel production⁶². Each Excel file contains the worksheet *alldata_flat*, in which all
 395 data points have been converted into a row-by-row format and can therefore be easily processed
 396 by machines.

397 In addition, in the application cases, both workflows use PV and wind generation time series
 398 from Pfenninger and Staffell⁶³. The time series are geographically and temporally detailed at
 399 country level, in this case for Germany, and in an hourly resolution for all weather years between
 400 1980 and 2019. The data is provided by the website *renewables.ninja*⁶⁴ under the open CC
 401 BY-NC 4.0 license as CSV files, which are structured into two columns containing the timestamp
 402 and the normalized generation.

403 Figure 4 shows an abridged DataDesc document describing the data model of the PV gener-
 404 ation time series for Germany. The CSV file comprises a tabular structure consisting of the two
 405 columns *national* (lines 13-18) and *time* (lines 23-25). While *national* holds the capacity factor
 406 float values for Germany in the range between 0 and 1, *time* represents the temporal dimension
 407 indexing the values using the coordinated universal time (UTC) format.

408 Shared Transformation Functions

409 In both workflows the same fundamental micro-transformations of data need to be performed.
 410 These operations were implemented to be compatible with ioProc and the general ones were
 411 applied both in the REMix and ETHOS.FINE workflows. Both workflows share the following
 412 transformation steps (we denote the exact function name of the ioProc action in brackets for

```

1 {
2   "dataDescVersion": "0.1.0",
3   "openapi": "3.0.0",
4   ...
5   "dataSchema" :
6   {
7     "type" : "object",
8     "format" : {
9       "example" : "text/csv"
10    },
11    "properties" : [
12      {
13        "identifier" : "national",
14        "type" : "number",
15        "format" : "float",
16        "minimum" : 0.0,
17        "maximum" : 1.0,
18        "dimensions" : [ "time" ]
19      }
20    ],
21    "dimensions" : [
22      {
23        "identifier" : "time",
24        "type" : "string",
25        "format" : "YYYY-MM-DD HH:MM:SS "
26      }
27    ]
28  }
29 }

```

Figure 4: Abridged DataDesc document describing the data model of a country's PV generation time series in an hourly resolution as provided by *renewables.ninja*⁶⁴.

413 each step):

414 The workflow starts with reading the source data into a memory data structure for further
415 processing (action *read_excel()*). We then clean the input data from whitespaces and special
416 characters in the text based raw data (action *strip_data()*) and convert numerical data to python
417 numerical types (action *remove_non_digit_values()*). The next steps are the reduction of the
418 dataset to the technologies we are interested in (in this case gas turbines, Li-ion batteries, heat
419 pumps and electrolyzers) (action *filter_technologies()*), to the relevant technology parameters
420 (action *filter_parameters()*) and finally to the target years (action *filter_years()*) and the mean
421 estimate (action *filter_estimate()*). Finally we exclude all columns with irrelevant data for our use
422 case from the dataset (action *drop_columns()*).

423 In the next phase of the data preparation, we have to make adjustments of the dataset la-
424 bels for example to be compatible with the REMix conventions. We hence rename the tech-
425 nologies (action *rename_technologies()*) and parameters (action *rename_parameters()*) to their
426 given standard names. Afterwards, we convert all values to consistent base units (action *con-
427 vert_units()*) and discount the costs to a uniform reference year (action *convert_currency()*). For
428 the cost conversion we delegate the task to the Currency Conversion for Python (CuCoPy) li-
429 brary^{52,65}, which we call from inside of the action. CuCoPy is an open-source software package
430 developed to streamline the integration of heterogenous currency data into research software.
431 The package enables the conversion of currencies from over 260 countries, along with the ability
432 to adjust financial data for inflation dating back to 1960, on an annual basis. As CuCoPy is inde-
433 pendently released under the open MIT license, it seamlessly integrates into existing applications

434 but also works autonomously within processing workflows like the ones described here.

435 **REMix Workflow Utilizing ioProc**

436 REMix is a modular codebase that allows to set up linear optimization models and was developed
437 by the Institute of Networked Energy Systems (VE) at the German Aerospace Centre (DLR) to
438 address research questions in the field of energy systems modeling.⁵⁶ The framework is pro-
439 grammed in Python and relies on GAMS as its mathematical programming language. REMix
440 can be used to address problems related to the optimal design and operation of future integrated
441 energy systems, both for target systems and transformation pathways. It implements methodolo-
442 gies like Pareto optimal search, modeling to generate alternatives, perfect foresight and myopic
443 optimization. One of the most important strengths of REMix is its high-performance scalability,
444 since it is designed for processing large spatio-temporal optimization scopes. This makes it fea-
445 sible to conduct path optimizations of high complexity with multiple years in hourly resolution and
446 a high spatial resolution including balance area sizes of up to continental level and individual
447 transmission lines. REMix is based on an implementation agnostic data model with makes it
448 applicable for different scopes. The data model consists of five basic building blocks, nodes and
449 goods. Nodes represent installations or connection points and goods describe everything that is
450 exchanged withing the system. Converters balance the consumption and production of goods,
451 storages allow the retention of goods across time, transfer links allow the transfer of these goods
452 between nodes. Sources and sinks allow respectively the input and output of goods from and
453 into the modeling scope. Complemented by indicators, users can model complex accounting
454 relations like carbon taxes or electricity market prices.

455 The primary input to REMix models consists of structured data, stored in CSV or DAT files,
456 which are loaded into the REMix model at the beginning of each model execution. These
457 datasets are comprehensive, covering all aspects of the energy systems model being analyzed.
458 They also play a crucial role in the configuration of the model, as REMix is a data-driven frame-
459 work. Therefore, the combination of a specific REMix software version and a complete dataset
460 constitutes a REMix instance, which is reproducible. In most cases, scientists do not create
461 entirely new datasets from scratch, since research questions rarely focus on completely novel or
462 previously not modeled systems. Instead, most modeling work is derived from a basic energy
463 system model that defines the scope and area of interest. This requires scientists to provide
464 data for specific parts of the model, such as particular technologies or regions, and properly link
465 them to the overarching energy system data. Technically, this often involves modifying or replac-
466 ing certain parts of an existing dataset. The core dataset that represents the energy system is
467 referred to as the baseline dataset. For a dataset to be usable as a baseline, it must conform
468 to the REMix conventions and be complete enough to execute the model without the need for
469 additional data.

470 This application case focuses on the common task of adding specific data to a baseline model
471 in REMix. To this end, we developed an exemplary workflow that integrates publicly available
472 open data into a baseline dataset. The main tasks involved are transforming the data to fit REMix
473 conventions and creating a new set of REMix-compatible input data files. We implemented this
474 workflow in ioProc, which consists of a series of distinct, specific actions. The implementation is
475 independent of any particular workflow or model context, enabling the exchange and sharing of
476 source code between the REMix and ETHOS.FINE workflows.

477 The REMix workflow begins with 12 general data cleaning and transformation actions that
478 focus on parsing and interpreting the input data format. These general actions are followed by
479 REMix-specific steps to clean and transform the raw input data to comply with REMix conven-
480 tions. In the final stage, the data is converted into the standard input data structure for REMix:
481 a two-dimensional table format with standardized labels and columns. The next step involves

482 incorporating the new dataset into the baseline model. To do this, we read the baseline dataset,
483 parameterize the entire model instance in REMix, and replace the relevant subset of the existing
484 baseline data with the new, more detailed dataset. In this example, the new dataset introduces
485 additional technologies, which must also be declared in REMix. This step completes the data
486 preparation process. At this point, the first dataset is ready to be written as REMix input files,
487 marking the final step of the workflow.

488 In this example, we incorporate two datasets, with the second being timeseries data. The
489 timeseries data is handled using the same workflow described above, though the required trans-
490 formations are fewer since the data is already in a format compatible with REMix. As a result,
491 we can skip the cleaning and conversion steps. The next steps involve reducing the dataset to
492 the target year for our model run and converting the labels to the REMix standard. The data is
493 now ready to be integrated into the base REMix model, following the same procedure as with
494 the technology dataset. We have now created a complete REMix input dataset consisting of the
495 baseline dataset, along with the incorporated technology and timeseries data. These files are
496 now ready for processing by REMix.

497 Reading a dataset into REMix is facilitated by a call to the *read_remix_csv* function, which we
498 have documented as part of this project with the DataDesc schema to provide additional meta
499 information and make it machine readable. The description includes a detailed description of the
500 two input parameters *file* and *schema* which were described as separate variables with their own
501 data schemas. While *file* is a simple string that refers to a CSV file *schema* is a complex data
502 structure which refers to a JSON file that is composed of various fields (including *name*, *title*,
503 *type*, *isAbout*, etc.). The return value of the function is the REMix interface internal data format,
504 which is a *pandas.DataFrame*.

505 To demonstrate how specific operations in one workflow can be generalized, such that they
506 become applicable in another workflow, we implemented two operations in the action format of
507 ioProc and used them in also in the ETHOS.FINE workflow. The selected actions are the reading
508 of the raw source data into memory (more precisely into pandas DataFrames) and the currency
509 conversion, which are based on CuCoPy. This demonstrates, that atomic transformation op-
510 erations, can be used across different workflow implementations and can be used to create a
511 shared basis of operations.

512 The full ioProc workflow encompasses the raw input datasets and the baseline REMix model
513 data, the workflow specification in the ioProc YAML file format, and the action folder containing all
514 ioProc actions. This bundle is published alongside with this publication on GitHub and constitutes
515 a reproducible workflow artefact.⁶⁶ With the documented ioProc version, it is now possible to
516 reproduce the workflow. Furthermore, in the spirit of good scientific practice, the action folder,
517 containing all operations needed for the workflow, is included and describes all scientific relevant
518 transformation of the data. It therefore constitutes a different form of documentation, which can
519 be used to examine the process in the future independent of an ioProc installation.

520 This approach supports reproducibility and transparency, as a separation of the basic data,
521 the workflow software components and the scientifically relevant transformations is achieved.
522 The ioProc actions can be used in the future to rebuild the workflow in any python context since
523 the ioProc actions do not depend on an ioProc installation. But also the reverse is true. As
524 easy as it is to use an ioProc action in a script, it is to modify existing source code in notebooks
525 or scripts to be usable as ioProc actions. This enables, without compromising the ability to
526 continue using the same code in existing scripts. The transformation into ioProc actions comes
527 with the added benefit that the meta data handling, logging and limited change tracking of ioProc
528 workflows are applicable without additional work.

529 From a research software engineering perspective, the ability to use ioProc action in other
530 environments is an enabling feature, as the individual testing of actions becomes easy to im-
531 plement. Such tests also serve as another form of scientific documentation beyond their usual

532 benefit of guaranteeing technical correctness. ioProc provides to scientist a framework in which
533 to build flexible and granular transformation operations, which can be combined and extended
534 as needed and applied in a broad range of software environments reaching from scripts and
535 notebooks to full software applications.

536 ETHOS.FINE Workflow Utilizing Jupyter Notebook

537 The *Framework for Integrated Energy Assessment* is part of the *Energy Transformation Path-*
538 *way Optimization Suite* at the division Jülich Systems Analysis of the the Institute of Energy
539 and Climate Research at Forschungszentrum Jülich.⁶⁷ The Python-based framework serves as
540 the basis for a diverse set of energy systems analyses conducted at the institute.^{68–72} Their
541 outcomes support investment decisions and policy-making in the transformation towards fully re-
542 newable energy systems. ETHOS.FINE leverages the capabilities of the general framework for
543 linear optimization, Pyomo, to model linear optimization problems for energy systems analysis.
544 ETHOS.FINE introduces the concepts of source, sink, storage and transmission components at
545 different spatial and temporal resolutions as building blocks of an energy systems model. Typi-
546 cal problems that can be analyzed with such models are cost-optimized future energy systems
547 of single buildings⁶⁸, municipalities^{69,70}, countries^{71–73}, or world-wide energy carrier transport⁷⁴.
548 Furthermore, ETHOS.FINE incorporates methods for spatial and temporal aggregation⁵⁴, tack-
549 ling the problem of long calculation times for large-scale mixed-integer linear energy system
550 optimization problems.

551 The relevant model information is stored in the *EnergySystemModel* container class. This
552 class holds general information such as units and location lists, as well as instructions for the
553 solving algorithm. In addition, instances of components are added to the *EnergySystemModel*
554 class in a modular way with the help of the *EnergySystemModel.add()* method. The output of the
555 optimization consists of values for the optimal design and operation of a minimum-cost system.

556 For the documentation of the programmatic interface, the constructors of the *EnergySys-*
557 *temModel* class and the component classes were described using DataDesc (*FINE.json* in the
558 supporting material).⁶⁶ Lines 1-243 contain general metadata of the software. From line 264
559 onward the *EnergySystemModel.add()* constructor function is described. The data schema de-
560 scription includes an explanation of the properties, types, ranges, and default values and in-
561 dicates required properties. Based on an example from the ETHOS.FINE repository⁵⁷, two
562 Jupyter notebooks have been created that contain the processing of raw input data to the for-
563 mat that is needed for the model (*01_data_processing.ipynb*) as well as the model generation
564 (*02_model_calculation.ipynb* in the supporting material). The preprocessing step reads time-
565 series for renewable energy potentials and energy demands from existing CSV files as well as
566 techno-economic parameters from a JSON file. The renewable energy potential timeseries col-
567 lected from the *renewables.ninja* website can be inserted into the models input CSV without
568 transformation because they are a unitless values between 0 and 1. Parameter values for an
569 absorption heat pump have been collected from the dataset of the Danish Energy Agency. Here,
570 the preprocessing methods described in the previous sections have been used for reading and
571 transforming the data. The raw data has been read from Excel into dictionary format with meth-
572 ods described in section *Shared Transformation Functions*. Economic values are described as
573 EUR/MW in the input data while the model uses EUR/GW. The transformation has been con-
574 ducted with the *convert_units()* method. Also, all monetary values are transformed to the year
575 2022 using CuCoPy.

576 In contrast to REMix, where one of the core functions with very few parameters was de-
577 scribed, only a small function was described for ETHOS.FINE, the *EnergySystemModel.add()*
578 function. The *add()* function expects many more parameters, some of which have custom com-
579 plex data types that are not found in the standard Python package library. These user-defined

580 data type descriptions can be summarized in unique schemas and referenced multiple times in
581 the document to avoid redundancy.

582 A helper script (*tools/comparison/main.py*) was used to automatically identify transformation
583 needs between the API and the datasets. The script was called via the command line, where
584 a DataDesc interface description, an input file description as well as a file path for the output
585 file were specified. As software descriptions usually consist of more than just one interface,
586 it was also necessary to specify which function and which parameter of said function should
587 be compared with the input data. The comparison is always carried out from the perspective
588 of the interface; meaning that a comparison is only successful if the input data serves all the
589 required parameters of the interface. On execution, JSON-formatted report files were generated
590 (see example in Figure 5), which provide information about overall compatibility (line 2), but
591 also more detailed information about missing expected data (lines 9-10) and mismatching data
592 formats (lines 18-23).

```
1 {  
2   "match": false,  
3   "detail": {  
4     "dataSchema": {  
5       "properties": [  
6         {  
7           "identifier": "windSpeed",  
8           "unit": {  
9             "expected": "https://qudt.org/vocab/unit/KILOM-PER-HR",  
10            "received": null  
11          },  
12        }  
13      ],  
14      "dimensions": [  
15        {  
16          "identifier": "time",  
17          "type": {  
18            "expected": "integer",  
19            "received": "string"  
20          },  
21          "format": {  
22            "expected": null,  
23            "received": "YYYY-MM-DD HH:MM:SS"  
24          },  
25        }  
26      ]  
27    }  
28  },  
29 }
```

Figure 5: Abridged DataDesc report⁶⁶ contrasting and comparing individual data model components of the ETHOS.FINE software interface and the wind generation time series from Pfenninger and Staffell⁶³ as provided by *renewables.ninja*⁶⁴.

593 Discussion

594 The coupling of research software and data into comprehensive workflows is common procedure
595 in the computational sciences to help answer complex research questions. Thereby, workflows

596 are not only relevant as versatile and reusable software tools, but are also a central component
597 of scientific exchange when it comes to tracing, reproducing and interpreting results data and
598 ultimately increasing the reliability of derived findings.³

599 Leipzig et al.¹⁴ point out that metadata plays a special role in the context of reproducible com-
600 putational research. They characterize the description of input or raw data and file intermediates
601 as they are processed in workflows as an essential core task of metadata. However, since none
602 of the established metadata schemas provide a sufficient description of the data models used
603 in data sets and software interfaces, information that can be automatically processed to enable
604 easy coupling of components in workflows is rarely available. The presented refinement of the
605 DataDesc schema, which allows data models to be annotated not only in software interfaces
606 but also in data sets, closes this gap and thus represents a necessary addition to established
607 domain-agnostic metadata standards like Dublin Core or schema.org.

608 The added value resulting from the ability to formally describe information of this kind in the
609 form of metadata is manifested in the downstream applications that this enables. Leipzig et al.¹⁴
610 state file format and content sanity checks, which are defined by input metadata but implemented
611 at the workflow level, as a prominent use case that has not yet been realized. To address this
612 use case, the DataDesc framework was expanded to include a tool that not only automatically
613 compares data content and formats, but also data structures and value ranges, identifies dis-
614 crepancies as transformation requirements for successful integration, and describes them in a
615 machine-actionable form. The ioProc workflow manager presented here then offers the option
616 of addressing identified transformation requirements by bundling and integrating transformation
617 paths in the form of modular adapters and making these available for reuse. Based on the uti-
618 lized modular workflow concept for the integration of software and data components, DataDesc
619 and ioProc support researchers in their manual design of research software workflows.

620 A promising application based on this comes from the innovative field of computer-aided
621 workflow design and integration, which could use such machine-actionable information to auto-
622 matically identify suitable data conversion and integration paths for model chains and suggest
623 them to researchers. The web-based scientific data processing platform Galaxy,⁷⁵ for example,
624 serves as a platform for describing, sharing, and publishing scientific calculation processes and
625 is designed to facilitate their discoverability and reusability in an accessible and interoperable
626 manner. Kumar et al.⁷⁶ have developed an approach that enables the platform to suggest tool
627 combinations based on usage patterns identified using deep learning. At this point, DataDesc's
628 automated comparison of the data models would provide information about transformation re-
629 quirements and a qualitative assessment of the proposed workflows. Computational sciences
630 have access to huge amounts of heterogeneous research data and research software. These
631 can no longer be fully evaluated using conventional data and tool searches, which often leads
632 to low reuse rates and redundant developments.^{77,78} Against this background, the computer-
633 assisted selection and integration of components into one's own research workflows is of partic-
634 ular importance.

635 In addition to designing and implementing scientific software workflows, ensuring that they
636 are used transparently and comprehensibly in studies and analyses remains a key challenge. In
637 this context, the problem of reproducibility in research has already been comprehensively docu-
638 mented.⁷⁹ Leipzig et al.¹⁴ describe computational process pipelines as a widely used method for
639 performing scientific analyses, that encourages parameterization and configuration that promote
640 reproducibility. With ioProc, a lightweight workflow manager has been released that natively im-
641 plements the modular workflow concept for integrating research data and software in the form of
642 adapters within reproducible and reusable data pipelines. Overall, the concept and tool devel-
643 opments presented here aim to help close the gap between abstract guidelines, such as those
644 systematized in the FAIR principles, and the daily work reality of researchers and to balance the
645 effort of metadata documentation with the corresponding practical benefits in handling scientific

646 workflows.

647 **Limitations**

648 In this work, we have for the first time presented application examples to demonstrate in a coher-
649 ent and comprehensive manner how metadata-based design and modular coupling processes
650 interact and can be implemented. Although the FAIR principles and open source and open data
651 are already becoming increasingly widespread, it will take some time before the publication of
652 entire workflows and modular transformation functions becomes standard practice in science.
653 As the approach presented here benefits from the high availability of existing and well-annotated
654 software and data components to quickly identify and integrate needed components into the
655 user's own workflows, its effect is initially limited and will only develop its full potential with the
656 increasing expansion and use of digital research infrastructures such as workflow repositories
657 and software registries. The DataDesc schema and the ioProc workflow manager are comple-
658 mentary components in this context, which, through further integration into the infrastructure, will
659 offer even better embedded and holistic solutions for everyday use.

660 The approaches and tools presented in this paper aim at loosely coupling models via transfor-
661 mation adapters, with the goal of achieving flexible interoperability between calculation steps that
662 may differ in data formats, data types, syntax, systems, and execution times. No adaptation of
663 data sets or software to predefined data and interface standards is required. In application con-
664 texts where workflows are mostly hard coupled, the direct benefit of the presented approach is
665 limited. In co-simulation, for example, models are usually provided with standardized interfaces
666 that do not require further transformation when interconnected. Only when hard and loose cou-
667 pling of models is combined within hierarchical workflows should clearly annotated and reusable
668 adapters be used again.

669 **Outlook**

670 Based on this work, future studies could further develop both the presented implementation
671 and design aspects. A reasonable next step would be the automated annotation of data files,
672 starting with a selection of especially compatible formats and working from there. Also, further
673 metadata publication pipelines should be included in the framework in order to supplement ex-
674 isting connections to software publication platforms with those for data publications, such as
675 the Open Energy Databus⁸⁰. Correspondingly the creation of a public transformation function
676 repository or catalogue with a stable API together would greatly improve transparency and con-
677 tribute to findability and reusability, i.e., core values of the FAIR principles. Based on such a data
678 and transformation function availability, tools can be developed which automatically identify the
679 needed transformation functions for connecting two different data formats (such as the format
680 of a data source and the format of a model data input) and the ability to create lists of needed
681 and missing transformations for scientists to build their workflows from. This would then open
682 up the possibility to further investigate limitations and best practices for automated generation
683 of workflows in different scientific contexts. Additional tooling can then be developed like rec-
684 ommendation systems for transformation functions or alternative workflow formulations based
685 on external factors like license compatibility. A different direction for future work would be the
686 improvement of the workflow tools themselves, such that they include automatic tracing of data
687 modifications along a specified workflow and including this information into metadata formats.
688 This can then be expanded towards systems that trace also the reason behind modifications and
689 support scientists in adding contextual information during their workflows in a user-friendly way.

690 **Resource Availability**

691 **Lead Contact**

692 Further information and requests for resources on the DataDesc framework should be directed
693 to and will be fulfilled by Patrick Kuckertz (p.kuckertz@fz-juelich.de).

694 Further information and requests for resources on the ioProc workflow manager should be
695 directed to and will be fulfilled by Benjamin Fuchs
696 (benjamin.fuchs@dlr.de).

697 **Materials Availability**

698 This study did not generate new unique reagents.

699 **Data and Code Availability**

700 The DataDesc metadata schema and the source code of the developed comparison tool for
701 DataDesc documents are publicly available under the open MIT license at <https://github.com/FZJ-IEK3-VSA/DataDesc>.⁵¹ In addition, these resources have been made available in their
702 current version 1.0 in the JülichDATA repository under the CC0 public domain dedication at
703 <https://doi.org/10.26165/JUELICH-DATA/DLCYV5>.

704 The ioProc workflow manager in version 2.2.0 under MIT license was used for this work and
705 is publicly available at <https://pypi.org/project/ioproc/> and as a git repository at <https://gitlab.com/dlr-ve/esy/ioproc>. Furthermore the addon ioprocmeta was also used in this
706 work and is available under BSD-3-clause license at <https://pypi.org/project/ioprocmeta/>
707 and as a git repository at <https://gitlab.com/dlr-ve/esy/ioprocmeta>.

708 The example and documentation files related to the presented application case are published
709 in a git repository at <https://github.com/dlr-ve-esy/modelcouplingworkflows> under BSD-3-
710 clause and CC-BY 4.0 licenses.
711
712

713 **Acknowledgements**

714 The authors would like to thank the Federal Ministry for Economic Affairs and Climate Action of
715 Germany (BMWK) for supporting this work with a grant for the project LOD-GEOSS (03EI1005A-
716 G). Furthermore, the authors are grateful to the German federal government, the German state
717 governments, and the Joint Science Conference (GWK) for their funding and support as part of
718 the NFDI4Ing consortium, managed by the German Research Foundation (DFG) – 442146713.
719 This work was also supported by the Helmholtz Association as part of the program "Energy
720 System Design".

721 **Author Contributions**

722 Conceptualization: P.K., and B.F.; methodology: P.K., B.F., and J.S.; software: B.F., H.G., K.K.,
723 and J.S.; validation: K.K., H.G., J.S., B.F., and P.K.; investigation: P.K., B.F., J.S., and J.G.; data
724 curation: H.G., B.F., E.A.R., and J.S.; writing – original draft: P.K., B.F., J.G., E.A.R., K.K., H.G.,
725 and J.S.; writing - review & editing: P.K., B.F., J.S., H.G., K.K., E.A.R., J.G., H.C.G., J.M.W., P.J.,

726 and J.L.; visualization: P.K.; supervision: P.K., B.F., J.M.W., J.L., and P.J.; project administration:
727 P.K., B.F., J.L., and P.J.; funding acquisition: P.J., and D.S.

728 **Declaration of Interests**

729 The authors declare no competing interests.

730 **Declaration of generative AI and AI-assisted technologies**

731 During the preparation of this work the authors used the tools DeepL and ChatGPT to check
732 grammar and spelling in a few places, and to make minor improvements to readability and style.
733 After using these tools, the authors reviewed and edited the content as needed, and take full
734 responsibility for the content of the publication.

735 **References**

- 736 1. Askeland, M., Morch, A., Papadimitriou, C., Di Somma, M., Coccia, A., Pinel, D., Richard-
737 son, P., and Sforza, G. (2023). Workflow-based architecture for optimal planning of inte-
738 grated local multi-energy systems. In: 2023 International Conference on Smart Energy
739 Systems and Technologies (SEST). IEEE (1–6).
- 740 2. Liu, J., Braun, E., Döpmeier, C., Kuckertz, P., Ryberg, D. S., Robinius, M., Stolten, D., and
741 Hagemeyer, V. (2019). Architectural concept and evaluation of a framework for the efficient
742 automation of computational scientific workflows: An energy systems analysis example.
743 *Applied Sciences* 9, 728.
- 744 3. Pelsler, T., Weinand, J. M., Kuckertz, P., and Stolten, D. (2025). ETHOS.REFLOW: An open-
745 source workflow for reproducible renewable energy potential assessments. *Patterns*.
- 746 4. Pueblas, R., Kuckertz, P., Weinand, J. M., Kotzur, L., and Stolten, D. (2023).
747 ETHOS.PASSION: An open-source workflow for rooftop photovoltaic potential assessments
748 from satellite imagery. *Solar Energy* 265, 112094.
- 749 5. Pfenninger, S., Hirth, L., Schlecht, I., Schmid, E., Wiese, F., Brown, T., Davis, C., Gidden,
750 M., Heinrichs, H., Heuberger, C. et al. (2018). Opening the black box of energy modelling:
751 Strategies and lessons learned. *Energy Strategy Reviews* 19, 63–71.
- 752 6. Kieser, S., Brown, J., Zdobnov, E. M., Trajkovski, M., and McCue, L. A. (2020).
753 Atlas: a snakemake workflow for assembly, annotation, and genomic binning of
754 metagenome sequence data. *BMC Bioinformatics* 21, 257. [https://doi.org/10.1186/
755 s12859-020-03585-4](https://doi.org/10.1186/s12859-020-03585-4).
- 756 7. Mölder, F., Jablonski, K. P., Letcher, B., Hall, M. B., Tomkins-Tinch, C. H., Sochat, V., Forster,
757 J., Lee, S., Twardziok, S. O., Kanitz, A., Wilm, A., Holtgrewe, M., Rahmann, S., Nahnsen,
758 S., and Köster, J. (2021). Sustainable data analysis with snakemake. *F1000Research* 10.
- 759 8. Pfenninger, S. (2017). Energy scientists must show their workings. *Nature* 542, 393–393.

- 760 9. Pelser, T., Weinand, J. M., Kuckertz, P., McKenna, R., Linssen, J., and Stolten, D. (2023).
761 Reviewing accuracy & reproducibility of large-scale wind resource assessments. *Advances*
762 *in Applied Energy* (100158).
- 763 10. Goble, C., Cohen-Boulakia, S., Soiland-Reyes, S., Garijo, D., Gil, Y., Crusoe, M. R., Peters,
764 K., and Schober, D. (2020). Fair computational workflows. *Data Intelligence* 2, 108–121.
- 765 11. Pimentel, J. F., Murta, L., Braganholo, V., and Freire, J. (2019). A large-scale study about
766 quality and reproducibility of jupyter notebooks. In: 2019 IEEE/ACM 16th International
767 Conference on Mining Software Repositories (MSR). (507–517). [https://doi.org/10.](https://doi.org/10.1109/MSR.2019.00077)
768 [1109/MSR.2019.00077](https://doi.org/10.1109/MSR.2019.00077).
- 769 12. Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A.,
770 Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., Bouwman, J., Brookes,
771 A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T., Finkers, R.,
772 Gonzalez-Beltran, A., Gray, A. J. G., Groth, P., Goble, C., Grethe, J. S., Heringa, J., 't
773 Hoen, P. A. C., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S. J., Martone, M. E., Mons,
774 A., Packer, A. L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone, S.-A.,
775 Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M. A., Thompson, M., van der
776 Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K.,
777 Zhao, J., and Mons, B. (2016). The fair guiding principles for scientific data management
778 and stewardship. *Scientific Data* 3, 160018. <https://doi.org/10.1038/sdata.2016.18>.
- 779 13. Booshehri, M., Emele, L., Flügel, S., Förster, H., Frey, J., Frey, U., Glauer, M., Hastings,
780 J., Hofmann, C., Hoyer-Klick, C., Hülk, L., Kleinau, A., Knosala, K., Kotzur, L., Kuckertz,
781 P., Mossakowski, T., Muschner, C., Neuhaus, F., Pehl, M., Robinius, M., Sehn, V., and
782 Stappel, M. (2021). Introducing the open energy ontology: Enhancing data interpretation
783 and interfacing in energy systems analysis. *Energy and AI* 5, 100074. [https://doi.org/](https://doi.org/10.1016/j.egyai.2021.100074)
784 [10.1016/j.egyai.2021.100074](https://doi.org/10.1016/j.egyai.2021.100074).
- 785 14. Leipzig, J., Nüst, D., Hoyt, C. T., Ram, K., and Greenberg, J. (2021). The role of metadata
786 in reproducible computational research. *Patterns* 2.
- 787 15. Fuchs, B., Riehm, J., Nitsch, F., and Wulff, N. (2020). ioproc - a light-weight workflow
788 manager in python. <https://gitlab.com/dlr-ve/esy/ioproc>.
- 789 16. Chan, L. M., and Zeng, M. L. (2006). Metadata interoperability and standardization - A
790 study of methodology, part I: achieving interoperability at the schema level. *D Lib Mag.* 12.
791 <https://doi.org/10.1045/JUNE2006-CHAN>.
- 792 17. Dublin Core Metadata Initiative (DCMI). DCMI Metadata Terms.
793 <https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>.
- 794 18. DataCite – International Data Citation Initiative e.V.. DataCite Metadata Schema 4.5.
795 <https://schema.datacite.org/>.
- 796 19. World Wide Web Consortium (W3C). Data Catalog Vocabulary (DCAT) - Version 3.
797 <https://www.w3.org/TR/vocab-dcat-3/>.
- 798 20. Schema.org. Schema.org Dataset. <https://schema.org/Dataset>.
- 799 21. Soiland-Reyes, S., Sefton, P., Crosas, M., Castro, L. J., Coppens, F., Fernández, J. M.,
800 Garijo, D., Grüning, B. A., Rosa, M. L., Leo, S., Carragáin, E. Ó., Portier, M., Trisovic,
801 A., Community, R., Groth, P., and Goble, C. A. (2022). Packaging research artefacts with
802 ro-crate. *Data Sci.* 5, 97–138. <https://doi.org/10.3233/DS-210053>.

- 803 22. Frictionless Data. Data Package (v1). <https://specs.frictionlessdata.io/>.
- 804 23. Hülk, L., Huber, J., Hofmann, C., and Muschner, C. (2024). Open Energy Family - Open
805 Energy Metadata (OEMetadata). GitHub. [https://github.com/OpenEnergyPlatform/](https://github.com/OpenEnergyPlatform/oemetadata)
806 [oemetadata](https://github.com/OpenEnergyPlatform/oemetadata).
- 807 24. World Wide Web Consortium (W3C). The RDF Data Cube Vocabulary.
808 <http://www.w3.org/TR/vocab-data-cube/>.
- 809 25. Vu, B., Pujara, J., and Knoblock, C. A. (2019). D-repr: A language for describing and
810 mapping diversely-structured data sources to rdf. In: Proceedings of the 10th International
811 Conference on Knowledge Capture. K-CAP '19 New York, NY, USA: Association for Com-
812 puting Machinery. ISBN 9781450370080 (189–196). [https://doi.org/10.1145/3360901.](https://doi.org/10.1145/3360901.3364449)
813 [3364449](https://doi.org/10.1145/3360901.3364449).
- 814 26. Garijo, D., Ratnakar, V., Gil, Y., and Khider, D.. The software description ontology. revision:
815 1.9.0. <https://w3id.org/okn/o/sd/1.9.0>.
- 816 27. The HDF Group. Hierarchical Data Format, version 5. <https://www.hdfgroup.org/HDF5/>.
- 817 28. Unidata. Network Common Data Form (NetCDF). <https://doi.org/10.5065/D6H70CW6>.
- 818 29. MLCommons. Croissant Format Specification - Version 1.0.
819 <http://mlcommons.org/croissant/1.0>.
- 820 30. Kuckertz, P., Göpfert, J., Karras, O., Neuroth, D., Schönau, J., Pueblas, R., Ferenz, S., En-
821 gel, F., Pflugradt, N., Weinand, J. M., Nieße, A., Auer, S., and Stolten, D. (2024). Datadesc:
822 A framework for creating and sharing technical metadata for research software interfaces.
823 *Patterns* 5. <https://doi.org/10.1016/j.patter.2024.101064>.
- 824 31. Foundation, F. S.. Gnu diffutils - comparing and merging files.
825 <https://www.gnu.org/software/diffutils/manual/>.
- 826 32. Project, G.. git-diff - show changes between commits, commit and working tree, etc.
827 <https://git-scm.com/docs/git-diff>.
- 828 33. Levenshtein, V. I. (1966). Binary Codes Capable of Correcting Deletions, Insertions and
829 Reversals. *Soviet Physics Doklady* 10, 707.
- 830 34. Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word rep-
831 resentations in vector space. Preprint at arXiv. [https://doi.org/10.48550/arXiv.1301.](https://doi.org/10.48550/arXiv.1301.3781)
832 [3781](https://doi.org/10.48550/arXiv.1301.3781).
- 833 35. Pennington, J., Socher, R., and Manning, C. D. (2014). GloVe: Global Vectors for Word
834 Representation. In: Proceedings of the 2014 conference on empirical methods in natural
835 language processing (EMNLP). (1532–1543).
- 836 36. Fellbaum, C. WordNet: An Electronic Lexical Database. The MIT Press (1998). ISBN
837 9780262272551. <https://doi.org/10.7551/mitpress/7287.001.0001>.
- 838 37. Möbius, M., Pependicker, K. J., and Böning, S. (2022). A-match. Zenodo. [https://doi.](https://doi.org/10.5281/zenodo.6641652)
839 [org/10.5281/zenodo.6641652](https://doi.org/10.5281/zenodo.6641652).
- 840 38. World Wide Web Consortium (W3C). Shapes Constraint Language (SHACL).
841 <https://www.w3.org/TR/shacl/>.

- 842 39. Gil, Y., Garijo, D., Khider, D., Knoblock, C. A., Ratnakar, V., Osorio, M., Vargas, H., Pham,
843 M., Pujara, J., Shbita, B., Vu, B., Chiang, Y.-Y., Feldman, D., Lin, Y., Song, H., Kumar,
844 V., Khandelwal, A., Steinbach, M., Tayal, K., Xu, S., Pierce, S. A., Pearson, L., Hardesty-
845 Lewis, D., Deelman, E., Silva, R. F. D., Mayani, R., Kemanian, A. R., Shi, Y., Leonard, L.,
846 Peckham, S., Stoica, M., Cobourn, K., Zhang, Z., Duffy, C., and Shu, L. (2021). Artificial
847 Intelligence for Modeling Complex Systems: Taming the Complexity of Expert Models to
848 Improve Decision Making. *ACM Transactions on Interactive Intelligent Systems* 11, 1–49.
849 <https://doi.org/10.1145/3453172>.
- 850 40. Yoo, A. B., Jette, M. A., and Grondona, M. (2003). Slurm: Simple linux utility for resource
851 management. In: Feitelson, D., Rudolph, L., and Schwiegelshohn, U., eds. *Job Scheduling*
852 *Strategies for Parallel Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN
853 978-3-540-39727-4 (44–60).
- 854 41. Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B., Bussonnier, M., Frederic, J., Kelley,
855 K., Hamrick, J., Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., Willing, C., and
856 development team, J. (2016). Jupyter notebooks - a publishing format for reproducible
857 computational workflows. In: Loizides, F., and Schmidt, B., eds. *Positioning and Power in*
858 *Academic Publishing: Players, Agents and Agendas*. Netherlands: IOS Press (87–90).
- 859 42. Köster, J., and Rahmann, S. (2012). Snakemake—a scalable bioinformatics workflow en-
860 gine. *Bioinformatics* 28, 2520–2522. <https://doi.org/10.1093/bioinformatics/bts480>.
- 861 43. Chang, M., Lund, H., Thellufsen, J. Z., and Østergaard, P. A. (2023). Perspectives on
862 purpose-driven coupling of energy system models. *Energy* 265, 126335. <https://doi.org/10.1016/j.energy.2022.126335>.
- 864 44. Wierling, A., Schwanitz, V. J., Altinci, S., Balazinska, M., Barber, M. J., Biresselioglu, M. E.,
865 Celino, M., Demir, M. H., Dennis, R., Dintzer, N., el Gammal, A., Fernandes-Peruchena,
866 C. M., Gilcrease, W., Gladysz, P., Hoyer-Klick, C., Joshi, K., Kruczek, M., Lacroix, D., Mol-
867 gorzata, M., Mayo-Garcia, R., Morrison, R., Paier, M., Perenato, G., Ramakrishan, M.,
868 Reid, J., Sciallo, A., Solark, B., Suna, D., Süß, W., Unger, A., Fernandez Vanoni, M. L., and
869 Vasiljevic, N. (2021). Fair metadata standards for low carbon energy research - a review of
870 practices and how to advance. *Energies* 14, 6692. <https://doi.org/10.3390/en14206692>.
- 871 45. Schriml, L. M., Chuvochina, M., Davies, N., Eloë-Fadrosch, E. A., Finn, R. D., Hugen-
872 holtz, P., Hunter, C. I., Hurwitz, B. L., Kyrpides, N. C., Meyer, F., Mizrachi, I. K., San-
873 sone, S.-A., Sutton, G., Tighe, S., and Walls, R. (2020). Covid-19 pandemic reveals the
874 peril of ignoring metadata standards. *Scientific Data* 7, 188. <https://doi.org/10.1038/s41597-020-0524-5>.
- 876 46. Forschungsgemeinschaft, D. (2022). Guidelines for Safeguarding Good Research Practice.
877 Code of Conduct. Deutsche Forschungsgemeinschaft. <https://doi.org/10.5281/zenodo.6472827>.
- 879 47. Federal Ministry of Education and Research. Aktionsplan Forschungsdaten.
880 [https://www.bmbf.de/bmbf/de/forschung/digitale-wirtschaft-und-gesellschaft/aktionsplan-](https://www.bmbf.de/bmbf/de/forschung/digitale-wirtschaft-und-gesellschaft/aktionsplan-forschungsdaten/aktionsplan-forschungsdaten_node.html)
881 [forschungsdaten/aktionsplan-forschungsdaten_node.html](https://www.bmbf.de/bmbf/de/forschung/digitale-wirtschaft-und-gesellschaft/aktionsplan-forschungsdaten/aktionsplan-forschungsdaten_node.html).
- 882 48. European Commission. Open science. https://rea.ec.europa.eu/open-science_en.
- 883 49. Kuckertz, P., Göpfert, J., Pueblas, R., Schönau, J., Weinand, J. M., and Stolten, D. (2023).
884 Software and data integration. Forschungszentrum Jülich GmbH. <https://doi.org/10.5281/zenodo.8117988>.
- 885

- 886 50. Booshehri, M., Emele, L., Flügel, S., Förster, H., Frey, J., Frey, U., Glauer, M., Hastings,
887 J., Hofmann, C., Hoyer-Klick, C. et al. (2021). Introducing the open energy ontology: En-
888 hancing data interpretation and interfacing in energy systems analysis. *Energy and AI* 5,
889 100074.
- 890 51. Kuckertz, P., Göpfert, J., Karras, O., Neuroth, D., Schönau, J., Pueblas, R., Ferenz, S., En-
891 gel, F., Pflugradt, N., Weinand, J. M., Kotzur, L., Nieße, A., Auer, S., and Stolten, D. (2024).
892 DataDesc - A framework for machine-actionable software metadata. Forschungszentrum
893 Jülich GmbH. <https://github.com/FZJ-IEK3-VSA/DataDesc>.
- 894 52. Schönau, J., Kuckertz, P., Weinand, J. M., Kotzur, L., and Stolten, D. (2024). Currency
895 Conversion for Python (CuCoPy). Forschungszentrum Jülich GmbH. <https://github.com/FZJ-IEK3-VSA/CuCoPy>.
- 897 53. Patil, S., Kotzur, L., and Stolten, D. (2022). Advanced spatial and technological aggrega-
898 tion scheme for energy system models. *Energies* 15, 9517. <https://doi.org/10.3390/en15249517>.
- 900 54. Hoffmann, M., Kotzur, L., and Stolten, D. (2022). The pareto-optimal temporal aggregation
901 of energy system models. *Applied Energy* 315, 119029.
- 902 55. Forschungszentrum Jülich GmbH - Jülich Systems Analysis (ICE-2) (2022). tsam - Time
903 Series Aggregation Module. Forschungszentrum Jülich GmbH. <https://github.com/FZJ-IEK3-VSA/tsam>.
- 905 56. German Aerospace Center (DLR) - Institute of Networked Energy Systems (2022). REMix -
906 Framework for energy system optimization models. GitLab. <https://gitlab.com/dlr-ve/esy/remix/framework>.
- 908 57. Forschungszentrum Jülich GmbH - Jülich Systems Analysis (ICE-2) (2023). ETHOS.FINE
909 - Framework for Integrated Energy System Assessment. Forschungszentrum Jülich GmbH.
910 <https://github.com/FZJ-IEK3-VSA/FINE>.
- 911 58. Forschungszentrum Jülich GmbH - Jülich Systems Analysis (ICE-2). ReadTheDocs Doc-
912 umentation of ETHOS.FINE - Framework for Integrated Energy System Assessment.
913 <https://vsa-fine.readthedocs.io/en/master/>.
- 914 59. Danish Energy Agency (2025). The danish energy agency. Danish Energy Agency. <https://ens.dk/en>.
- 916 60. Danish Energy Agency (2024). Technology data - energy plants for electricity and district
917 heating generation. Danish Energy Agency. <https://ens.dk/media/6378/download>.
- 918 61. Danish Energy Agency (2024). Technology data – energy storage. Danish Energy Agency.
919 <https://ens.dk/media/6446/download>.
- 920 62. Danish Energy Agency (2024). Technology data – renewable fuels. Danish Energy Agency.
921 <https://ens.dk/media/6444/download>.
- 922 63. Pfenninger, S., and Staffell, I. (2016). Long-term patterns of european pv output using 30
923 years of validated hourly reanalysis and satellite data. *Energy* 114, 1251–1265.
- 924 64. Pfenninger, Stefan and Staffell, Iain. Renewables.ninja. <https://www.renewables.ninja/>.

- 925 65. Weinand, J. M., Hoffmann, M., Göpfert, J., Terlouw, T., Schönau, J., Kuckertz, P., McKenna,
926 R., Kotzur, L., Linßen, J., and Stolten, D. (2023). Global Icoes of decentralized off-grid
927 renewable energy systems. *Renewable and Sustainable Energy Reviews* 183, 113478.
- 928 66. Kuckertz, P., Fuchs, B., Schönau, J., Gardian, H., Knosala, K., Arellano Ruiz, E., Göpfert,
929 J., Gils, H. C., Weinand, J. M., Jochem, P., Linßen, J., and Stolten, D. (2025). Supplemen-
930 tary material to 'Model coupling through reproducible adapter workflows based on shared
931 transformation functions'. Forschungszentrum Jülich GmbH, Deutsches Zentrum für Luft-
932 und Raumfahrt (DLR). <https://github.com/dlr-ve-esy/modelcouplingworkflows>.
- 933 67. Groß, T., Knosala, K., Hoffmann, M., Pflugradt, N., and Stolten, D. (2023). Ethos.fine: A
934 framework for integrated energy system assessment. Preprint at arXiv. [https://doi.org/
935 10.48550/arXiv.2311.05930](https://doi.org/10.48550/arXiv.2311.05930).
- 936 68. Omoyele, O., Matrone, S., Hoffmann, M., Ogliari, E., Weinand, J. M., Leva, S., and Stolten,
937 D. (2024). Impact of temporal resolution on the design and reliability of residential energy
938 systems. *Energy and Buildings* 319, 114411.
- 939 69. Risch, S., Weinand, J. M., Schulze, K., Vartak, S., Kleinebrahm, M., Pflugradt, N., Kullmann,
940 F., Kotzur, L., McKenna, R., and Stolten, D. (2024). Scaling energy system optimizations:
941 Techno-economic assessment of energy autonomy in 11 000 german municipalities. *Energy
942 Conversion and Management* 309, 118422.
- 943 70. Weinand, J. M., Vandenberg, G., Risch, S., Behrens, J., Pflugradt, N., Linßen, J., and
944 Stolten, D. (2023). Low-carbon lithium extraction makes deep geothermal plants cost-
945 competitive in future energy systems. *Advances in Applied Energy* 11, 100148.
- 946 71. Jacob, R., Hoffmann, M., Weinand, J. M., Linßen, J., Stolten, D., and Müller, M. (2023). The
947 future role of thermal energy storage in 100% renewable electricity systems. *Renewable
948 and Sustainable Energy Transition* 4, 100059.
- 949 72. Tsani, T., Pelsler, T., Ioannidis, R., Maier, R., Chen, R., Risch, S., Kullmann, F., McKenna,
950 R., Stolten, D., and Weinand, J. (2024). Out of sight, out of mind? cost of minimizing
951 visibility of nationwide renewable energy systems. Preprint at Research Square. [https:
952 //doi.org/10.21203/rs.3.rs-5017073/v1](https://doi.org/10.21203/rs.3.rs-5017073/v1).
- 953 73. Schöb, T., Kullmann, F., Linßen, J., and Stolten, D. (2023). The role of hydrogen for a
954 greenhouse gas-neutral germany by 2045. *international journal of hydrogen energy* 48,
955 39124–39137.
- 956 74. Franzmann, D., Heinrichs, H., Lippkau, F., Addanki, T., Winkler, C., Buchenberg, P.,
957 Hamacher, T., Blesl, M., Linßen, J., and Stolten, D. (2023). Green hydrogen cost-potentials
958 for global trade. *international journal of hydrogen energy* 48, 33062–33076.
- 959 75. Community, T. G. (2024). The galaxy platform for accessible, reproducible, and collaborative
960 data analyses: 2024 update. *Nucleic Acids Research* 52, W83–W94. [https://doi.org/10.
961 1093/nar/gkae410](https://doi.org/10.1093/nar/gkae410).
- 962 76. Kumar, A., Rasche, H., Grüning, B., and Backofen, R. (2021). Tool recommender system in
963 galaxy using deep learning. *GigaScience* 10, g1aa152.
- 964 77. Habermann, T. (2020). Metadata and reuse: Antidotes to information entropy. *Patterns* 1.

- 965 78. Palmblad, M., Lamprecht, A.-L., Ison, J., and Schwämmle, V. (2019). Automated workflow
966 composition in mass spectrometry-based proteomics. *Bioinformatics* 35, 656–664.
- 967 79. Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. Nature Publishing Group UK
968 London.
- 969 80. DBpedia Association affiliated with Institut für Angewandte Informatik e. V. (InfAI). Open
970 Energy Databus. <https://databus.openenergyplatform.org/>.