

# From Ontology to Metadata: A Crawler for Script-based Workflows

**HOMER: a tool for extraction and re-use of ontology-based metadata in high-performance measurement and computing workflows**

Giuseppe Chiapparino  <sup>1</sup>

Benjamin Farnbacher  <sup>1</sup>

Nils Hoppe  <sup>1</sup>

Radoslav Ralev  <sup>2</sup>

Vasiliki Sdralia  <sup>3</sup>

Christian Stemmer  <sup>1</sup>

1. TUM School of Engineering and Design; Department of Engineering Physics and Computation; Chair of Aerodynamics and Fluid Mechanics, Technical University of Munich, Garching; Germany.
2. TUM School of Computation, Information and Technology; Department of Informatics, Technical University of Munich, Garching; Germany.
3. TUM School of Engineering and Design; Department of Engineering Physics and Computation; Chair of Aerodynamics and Fluid Mechanics; Munich Data Science Institute (MDSI), Technical University of Munich, Garching; Germany.



**Date Received:**

2023-01-20

**Licenses:**

This article is licensed under: 

**Keywords:**

Metadata extraction, HPMC, Ontology, Research Data Management

**Data availability:**

Data can be found here: [https://gitlab.lrz.de/nfdi4ing/crawler/-/tree/master/SimpleApplication\\_PizzaOntology](https://gitlab.lrz.de/nfdi4ing/crawler/-/tree/master/SimpleApplication_PizzaOntology)

**Software availability:**

Software can be found here: [doi:10.14459/2022mp16944012](https://doi.org/10.14459/2022mp16944012)

**Abstract.** The present work introduces HOMER (**High Performance Measurement and Computing tool for Ontology-based Metadata Extraction and Re-use**), a python-written metadata crawler that allows to automatically retrieve relevant research metadata from script-based workflows on HPC systems. The tool offers a flexible approach to metadata collection, as the metadata scheme can be read out from an ontology file. Through minimal user input, the crawler can be adapted to the user’s needs and easily implemented within the workflow, enabling to retrieve relevant metadata. The obtained information can be further automatically post-processed. For example, strings may be trimmed by regular expressions or numerical values may be averaged. Currently, data can be collected from text-files and HDF5 files, as well as directly hardcoded by the user. However, the tool has been designed in a modular way, so that it allows straightforward extension of the supported file-types, the instruction processing routines and the post-processing operations.

## 1 Introduction

Nowadays, scientists are called to handle large amount of generated data, store them in repositories and distribute them among other scientists or the scientific community, something that makes it hard for them to keep track of all the relevant information over time and space. This can lead to the generation of a large quantity of forgotten and unused data, the so-called Dark Data [1], [2]. **Although every researcher implements some sort of Research Data Management (RDM), either consciously or unconsciously, to avoid the loss of precious information, standardized RDM**

8 approaches, such as the FAIR data principles (Findable, Accessible, Interoperable, and Re-usable),  
9 have been proposed in order to provide a more structured and potentially efficient solution to  
10 these problems. From the beginning of a project, the scientist should have a data-management  
11 plan on how the data will be organized, where they will be stored safely and who should be able  
12 to access the data and re-use them. In fact, accompanying the complete data-generation process  
13 with a proper data-management plan will have two benefits. On one side, it will be easier to  
14 reproduce old research works for future scientists. On the other side, well-documented data will  
15 enable effective secondary research.

16 For that reason, the German Federal Government funded NFDI (Nationale Forschungsdaten-  
17 infrastruktur [National Research Data Infrastructure]), to establish an infrastructure on RDM,  
18 providing an environment where scientists can develop solutions to research questions and make  
19 their findings and innovations sustainable by implementing the FAIR data principles. NFDI4Ing  
20 (NFDI für die Ingenieurwissenschaften [NFDI for Engineering] [3]), one of the consortia funded  
21 by the NFDI initiative, brings together the engineering communities to develop, standardise and  
22 provide methods and services to make engineering research data FAIR.

23 One major factor of making data FAIR is the implementation of a controlled vocabulary with com-  
24 mon terminology. The use of a controlled vocabulary is essential for findability, interoperability,  
25 and consequently, the re-use and the establishment of new user models. Most of the research  
26 data in the HPMC domain is neither documented nor are metadata sets available, as common  
27 terminologies for HPMC in the engineering sector still need to be developed and established  
28 within the community. HOMER allows to automatically retrieve relevant research metadata  
29 from script-based workflows on HPC systems and therefore supports researchers to collect  
30 and publish their research data within a controlled vocabulary using a standardized workflow.  
31 Controlled vocabularies, and the relations and restrictions between their terms, are practically  
32 implemented through the use of ontologies. An ontology defines a shared conceptualization  
33 of a common vocabulary, semantic relations of data and the syntactic as well as the semantic  
34 interoperability, including machine-interpretable definitions of basic concepts in the domain and  
35 the relations among them. The NFDI4Ing consortium has developed an ontology as a common  
36 classification of engineering data in a taxonomic hierarchy with standardized vocabulary and pro-  
37 cedures. Metadata4Ing (Metadata for Engineering [4]) aims at providing a thorough framework  
38 for the semantic description of research data, with a particular focus on engineering sciences  
39 and neighboring disciplines. Metadata4Ing re-uses elements from the existing terminologies  
40 and ontologies, such as DCMI Metadata Terms [5] or the PROV (Provenance Namespace)  
41 ontology [6]), whose terms were imported into Metadata4Ing. This ontology allows a thorough  
42 description of the whole data-generation process (experiment, observation, simulation), covering  
43 aspects such as: the object of investigation, all sample and data manipulation procedures, a  
44 summary of the data files and the information contained, and all personal and institutional roles.  
45 The NFDI4Ing framework entails many working groups called “archetypes”. Among them, the  
46 role of archetype DORIS is twofold: on one side, to create a HPMC-(sub)ontology based on  
47 Metadata4Ing in order to establish a consistent terminology for computational fluid-dynamics  
48 (CFD) workflows in high performance computing (HPC) systems; on the other side, to develop a  
49 metadata crawler, presented in this work, for metadata extraction. The expansion to an HPC-sub-  
50 ontology is based on modularity and fits in the primary Metadata4Ing classes of method, tool,

51 object of research. The expansion includes suggestions of unambiguous terms for domain-related  
52 metadata expressed in classes, object properties (relations) and data properties. These classes  
53 have been developed in a community-based approach and represent common methods and tools  
54 for workflows in engineering research on HPMC systems. The crawler named HOMER (HPMC  
55 tool for Ontology-based Metadata Extraction and Re-use) is intended as a RDM tool to automate  
56 the retrieval of metadata and is designed to be used in script-based HPMC applications.

57 In the field of RDM, many solutions and tools for metadata extraction have been proposed in  
58 the recent years. While all of them share with HOMER the same core concept of automating  
59 metadata extraction on HPC systems, they implement different approaches and solutions to the  
60 problem, introducing a great variety of capabilities.

61 For example, the RDM system at the University of Huddersfield, iCurate [7], provides a tailored  
62 solution to HPMC data with the functionalities of metadata retrieval, departmental archiving,  
63 workflow management system and metadata validation and self inferencing. This last functional-  
64 ity requires the metadata to be mapped onto a suitable ontology. iCurate offers support for all  
65 aspects of data management, but the actual extraction of metadata is limited to the annotations  
66 made by the user in a HPC job file. While this guarantees a non-intrusive integration within the  
67 workflow, it doesn't allow the user to retrieve metadata from output files.

68 Another tool for research-data management is represented by signac [8]. This is a lightweight  
69 framework providing all the components to create a searchable and shareable **dataspace (a**  
70 **decentralized infrastructure for data sharing and exchange based on commonly agreed principles**  
71 **[9])**. The core application is a semi-structured database that allows storing the original files on  
72 the file system along with the associated metadata, which is created on-the-fly and saved in  
73 human-readable format. The tool also allows for workflow management thanks to the signac-flow  
74 application. The framework is implemented in python and is designed to be used in HPC systems.  
75 However, in order to perform the metadata annotation, signac requires the user to wrap the  
76 original simulation code into a script.

77 The extraction of metadata from output files is possible with Xtract [10], [11]. In general, this  
78 powerful serverless middleware provides an effective, flexible and scalable way to retrieve  
79 metadata from **very large data lakes (centralized systems storing data in raw format [12])**. **With**  
80 **Xtract, metadata can be extracted both centrally ("central mode", i.e. fetching the metadata**  
81 **all at once from the different repositories where the data are generated) and at the edge ("edge**  
82 **mode", i.e. extracting and storing the metadata as soon as the data are created and at the location**  
83 **where they are generated)**. The service implements several different extractors (written in python  
84 or bash), which are run dynamically according to the type of file that needs to be crawled.  
85 This allows to handle a vast quantity of different data formats typically employed in scientific  
86 applications. Machine learning is used to infer the type of file to be crawled, so as to choose  
87 the best extractor(s) for that file in the shortest time possible. Finally, the tool is highly portable  
88 as it is wrapped in a docker container [13] and can be linked to Globus [14], [15]. Xtract is a  
89 powerful service, which is however more suited to data lakes, rather than to be applied within a  
90 workflow. Moreover, the information that can be retrieved is somewhat limited and not entirely  
91 customizable by the user.

92 From this point of view, more freedom is given by ExtractIng [16], a generic automated metadata-

93 extraction toolkit. Again suited for HPC systems, ExtractIng is a Java-written standalone tool  
94 that needs to be run once simulation outputs have been produced. It is easy to integrate within a  
95 workflow and offers both native and parallel implementation of the parsing algorithm, which  
96 makes the code scalable for HPC applications. The metadata extraction is based on the metadata  
97 scheme provided by EngMeta [17]. The tool is code-independent, in the sense that an external  
98 configuration file allows to adapt the metadata extraction to the specific simulation code and  
99 computational environment. While this provides a generic extraction tool, the configuration file  
100 needs to be manually written and adapted by the user (even though it only needs to be done once  
101 per code).

102 Another solution is represented by Brown Dog [18], which consists of two services called DAP  
103 and DTS. The first provides file conversion, while the second performs extraction and analysis  
104 of metadata. Brown Dog aims at leveraging already existing software, libraries and services in  
105 order to provide an automated aid in RDM. The implementation of an elasticity module provides  
106 for an optimized auto-scale of the two services based on the system demand. Moreover, a tool  
107 catalog points the user to the most suitable option for file conversion and metadata extraction.  
108 The extracted metadata is returned as a JSON file. However, this operation is performed by  
109 passing the data to the web-based service Clowder. This particular aspect might pose some  
110 limitations in the use of Brown Dog.

111 Some of the services that provide metadata extraction might be domain specific. For example,  
112 ScienceSearch [19] is a generalized and scalable search infrastructure, which employs machine  
113 learning to capture metadata. Information is retrieved not only from regular data, but also from  
114 the context and the surroundings artifacts (proposals, publications, file system structures and  
115 images) of the data, allowing for an enrichment of the extracted metadata. The service provides  
116 a web interface, where users can submit their text queries, and also provides the possibility  
117 for the users to give feedback on the collected metadata, so as to improve the search quality.  
118 However, the data model is unique to the NCEM dataset, which includes data relevant to the  
119 field of electron microscopy.

120 Within the NFDI environment, Swate [20] is an Excel add-in for the annotation of experimental  
121 data and computational workflows developed by the consortium NFDI4Plants. The tool is  
122 intended for metadata annotation based on the ontology provided by the user. The use of a  
123 spreadsheet environment aims at providing an intuitive and low-friction workflow, principally  
124 focusing on wet-lab applications, where the user can annotate work-relevant metadata as the  
125 experiment is performed. Hence, in this tool, the work is done manually by the user.

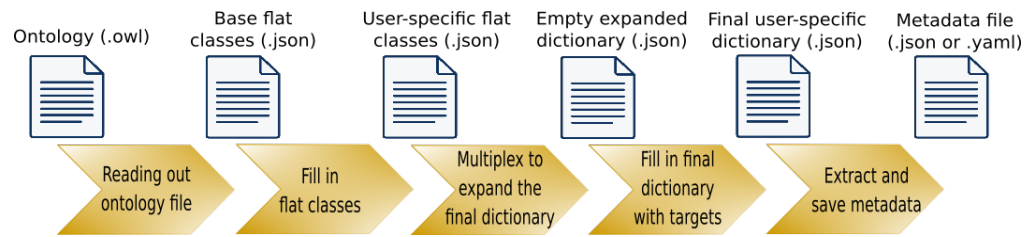
126 A second tool within the NFDI infrastructure is presented in this paper. Developed within the  
127 NFDI4Ing consortium at the Technical University of Munich (TUM), HOMER is a metadata  
128 crawler to be integrated in script-based (HPMC) workflows aiming at retrieving metadata that  
129 can be attached to the raw data published by researchers. **The tool is designed to be flexible and  
130 adjustable to the user's needs in its application and easy to implement in potentially any HPMC  
131 workflow. This development approach tries to overcome all the shortcomings highlighted for  
132 the RDM solutions and tool reviewed in this section, and to allow HOMER to be suitable for  
133 a wide range of applications. In fact, the crawler can retrieve metadata from text and binary  
134 (HDF5) files, as well as from user's annotations and terminal commands, at any stage of the**

135 workflow without interfering with the other processes composing the workflow. The automated  
136 extraction of metadata can be performed both in edge as well as in central mode, making the  
137 tool suitable for extracting information also from central repositories (such as data lakes). The  
138 metadata extraction is based on the ontology schemes chosen by the users. However, the users  
139 do not need to strictly adhere to a fixed scheme, but can adjust and customize it according to  
140 their needs. Moreover, although developed primarily keeping engineering sciences as the main  
141 use application, HOMER can be employed to retrieve metadata from HPMC workflows applied  
142 to a wide variety of research fields. Finally, the tool has been written with a modular structure,  
143 so it can easily be developed further to include new features. Hence, HOMER is proposed as a  
144 flexible and consistent RDM tool that can be used in a wide variety of applications and fields  
145 with limited user inputs in order to easily promote the FAIR principles and enrich the data created  
146 by the user.

147 The code structure is described in section 3, while a simple application is described in section 4.  
148 Finally, in section 5, an overview on the future steps in the code development is given.

## 149 **2 Characterization of the problem**

150 Many numerical applications, such as optimization problems or parametric studies, require the  
151 user to solve essentially the same problem with slightly different inputs each time. To make an  
152 example in the field of CFD, assume that it is necessary to assess the aerodynamic characteristics  
153 of an aircraft wing during different phases of the flight (take-off, climb, cruise and so on). Hence,  
154 the user will perform a certain number of simulations employing the same geometry of the  
155 wing while varying the freestream conditions (pressure, density, temperature, Reynolds number  
156 and so on) provided as input to the simulation. In such a case, especially when the number  
157 of simulations to be performed is large, automating the workflow (or parts of it) by means of  
158 script-based processes enables an efficient use of the available computing resources. For example,  
159 the user could create a script to change the freestream-input parameters as soon as a simulation  
160 ends so that the following one with new conditions would start immediately. Together with the  
161 data generation, the researcher should also aim at retrieving and storing the relevant metadata  
162 for all the computations performed, in order to comply with the FAIR principles and add value  
163 to the data gathered. In the wing-study example, the most obvious relevant metadata would be  
164 the different input freestream conditions associated to each simulation result, but the user could  
165 be also interested in storing information on the specific hardware or software (version of the  
166 code, version of the compiler, ID of the computational node, and so on) used for the simulations,  
167 for example. The information that needs to be extracted might be scattered across the different  
168 files that are usually generated during numerical calculations, such as the input and output files  
169 associated with the simulation, as well as files generated by the HPC system. Therefore, the user  
170 will have to go through all these files for each simulation and recursively extract and store the  
171 metadata. Doing such a job by hand would be certainly time consuming and would look as a  
172 viable option only if the number of simulations is very limited. In HPMC applications where  
173 hundreds of simulations are performed, this approach would be prohibitive to say the least and,  
174 therefore, the use of a dedicated extraction routine would be the preferential and most efficient  
175 choice. At this point, for the user it would be a matter of either writing an extraction routine



**Figure 1:** The five steps the crawler is currently composed of and related input/output files.

176 from scratch, which would guarantee a perfect compatibility even for very specific codes but  
 177 requires time and resources directly spent by the user, or employ an already available tool, at the  
 178 price of possible overheads to properly implement the tool within the workflow. In the latter  
 179 case, HOMER would come in handy as a valid support for the researchers. In fact, HOMER  
 180 is intended to be used, potentially, for any script-based application and is designed to be easily  
 181 adjusted to different research fields.

## 182 3 Code Description

### 183 3.1 Characteristics

184 HOMER is a code written in python and, at the moment, its complete extraction workflow  
 185 consists of 5 steps, as shown in figure 1. The code has been developed as a collection of modular  
 186 routines, each performing a different action, rather than as a single script. This guarantees a  
 187 flexible application of the crawler, as the user doesn't have always to perform all the steps  
 188 described in the next subsection, especially once the initial setup of the workflow has been  
 189 done for the first time. The modularity of the code also allows the developers and the users to  
 190 easily modify and expand the capabilities of the crawler, so that the code can be tailored for  
 191 specific applications, if needed. HOMER can be employed both locally and on HPC systems.  
 192 However, it should be noted that, as of now, it does not support a parallel implementation to  
 193 parse the target files. Moreover, the tool only covers the stages of planning, creating/collecting  
 194 and processing/analyzing within the data life cycle (figure 2).

195 Conceived as a tool to be integrated in script-based workflows, the crawler should be run after  
 196 the simulation (or, potentially, any processing step), similarly to ExtractIng. Hence, the metadata  
 197 are naturally extracted in edge mode (where the data are generated). However, the tool can also  
 198 be used to retrieve metadata from centrally-stored, previously collected data, similarly to Xtract.  
 199 In this case, though, the user has to perform some extra steps according to the specific case at  
 200 hand (an example is given at the end of section 4). Together with metadata extraction, the tool  
 201 gives the user the opportunity to perform some simple post-processing operations as well, such  
 202 as trimming strings or calculating the minimum, maximum or mean of a series of values.

### 203 3.2 Implementation

204 When running the code for the first time, five steps (figure 1) are needed, with two of them  
 205 requiring direct user input. The overview of this five-steps workflow is given in the next  
 206 paragraphs, while an application on a CFD-based example is shown in section 4.



**Figure 2: Data Life Cycle [21].**

207 The first step consists in reading the ontology file and creating an empty dictionary containing  
208 the flat classes as specified in the ontology. A “flat class” in this context is the initial state of  
209 the class in which the properties have not been specified, yet. The step is performed by the  
210 routine `ClassUtils.py` and takes advantage of the python package `Owlready2` to work on  
211 the ontology. The empty dictionary is a list of the classes and their attributes as they appear in  
212 the ontology file. The dictionary acts as a template, where all the flat classes are listed but not  
213 filled-in, yet. Hence, no specific instance of a class is created at this point. The file, however,  
214 contains the fields that allow the user to specify how many instances and related properties of  
215 each class are to be created.

216 The second step has to be performed manually and serves the purpose of preparing the dictionary  
217 file to be used by the “Multiplexer”, as explained in the third step. The user needs to specify  
218 how many instances of each class need to be retrieved. This is done by giving a numerical  
219 value to the keyword `__count__` in the corresponding class. Similarly, the user can specify how  
220 many properties each instance of a class needs to have by indicating the numerical value in the  
221 corresponding property list.

222 The third step consists in the Multiplexer and is performed by the routine `Multiplexer.py`.  
223 The output is the original flat-class dictionary which has been now expanded according to the  
224 needs of the user, so that the new file contains a list of all the needed instances for each class and  
225 all the properties for each instance of a class.

226 The fourth step again has to be performed manually. The user has to fill in the multiplexed  
227 dictionary by specifying, for each instance and property, where the crawler should look for  
228 the data and how it should retrieve them. This information is to be provided by specifying the

229 three keywords "path", "type" and "pattern". The filled expanded dictionary works as a  
230 configuration file for the final step.

231 The fifth and last step consists in the actual extraction of the metadata and is performed through  
232 the routine `EntityUtils.py`. Once the metadata have been extracted, it is printed out to a file  
233 in a structured human-readable format (JSON or YAML). Currently, metadata can be extracted  
234 from files, such as text files using regular expressions or HDF5 files using `h5py`, from the output  
235 of a operating-system command, or can be directly hardcoded during the fourth step, if needed.

236 As mentioned, the user must first configure the crawler to the specific simulation code in use  
237 (type and amount of metadata available to be extracted, location and format of the target files,  
238 extraction methods...). This involves the two (lengthy) manual steps just described. **However,**  
239 **once the first setup has been completed, the configuration file (created at the end of step four)**  
240 **can be re-used with little to no modification every time a new simulation is performed by the**  
241 **user and new metadata need to be extracted. Specifically, step five simply needs to be performed**  
242 **in the subsequent runs.** This allows for a seamless integration of the tool within the workflow.

## 243 **4 Example Application**

244 **In this section, a simple usage of the tool is shown for an application within the CFD field.**  
245 **Namely, the same test problem of a wing aerodynamic optimization mentioned in 2 is considered,**  
246 **in order to show the capabilities of the tool directly applied to an engineering application.**

247 **Nonetheless, the reader is also invited to try out the more generic step-by-step test case based on**  
248 **a simplified "Pizza ontology" available in the GitLab code repository (all the relevant files are**  
249 **in the directory `/SimpleApplication_PizzaOntology`). This purposely generic example is**  
250 **intended to provide a complete overview of the implementation of HOMER within a script-based**  
251 **workflow in a simple, clear and application-independent way. The ontology file used in the**  
252 **initial step is loosely based on the "Pizza ontology" provided by Stanford University in their**  
253 **Protégé tutorial [22].**

### 254 **4.1 Application to a CFD case**

255 **Taking the case of simulations on an airplane wing at different freestream conditions as an**  
256 **example, the user could be interested in extracting values such as freestream pressure, temperature,**  
257 **Mach and Reynolds number, as well as the ID-number of the node(s) where the calculation was**  
258 **performed and the software version used at the time of the calculation. This metadata information**  
259 **would then be attached to the results of the corresponding simulations to help making the data**  
260 **compliant with the FAIR principles.**

261 **In this example, the code NSMB [23] is employed to perform the simulations, and all the**  
262 **relevant pieces of information are extracted from the input and output files of the code and are**  
263 **classified according to a simplified version of the `Metadata4Ing` ontology for sake of simplicity.**  
264 **Therefore, in this example, the reference classes are limited to "Processing\_Step", "Tool"**  
265 **and "Method", with information such as parameters name and numerical value being considered**  
266 **as properties of those classes.**



267 In this example, it is assumed that the crawler is employed in edge mode, meaning that the  
 268 HOMER is invoked right after each simulation has finished to immediately extract the corre-  
 269 sponding metadata. This means that all the file paths specified by the keyword "path" are  
 270 relative paths, as the crawler runs from the same directory where the simulations are performed.  
 271 The example shows all the five steps described in section 3, which would be ideally required  
 272 for the first usage of the tool. After that, HOMER can be seamlessly integrated in the workflow  
 273 without further modifications.

274 After setting up the (optional) python virtual environment and having installed the crawler, the  
 275 first command to run is `ClassUtils.py`, which retrieves the ontology and creates a dictionary  
 276 with the flat classes (i.e. classes to be filled-out by the user and where, right after step 1, the  
 277 properties have simply placeholder values) in the form of a `.json` file. The content of such a  
 278 file would look like the one shown in the next lines.

```
279 {
280     "Processing_Step": {
281         "__count__": 1,
282         "__restrictions__": "",
283         "Name": [1],
284         "Parameter": [1],
285         "Has_numerical_value": [1],
286     },
287     "Tool": {
288         "__count__": 1,
289         "__restrictions__": "",
290         "System_component": [1],
291         "Name": [1],
292         "ID": [1],
293     },
294     ...
295 }
```

**Listing 1:** Example of the content inside the file produced after step 1.

296 This empty dictionary has to be manually adjusted to the specific case by the user. In this  
 297 example, only one processing step is foreseen (running the simulation), for which five parameters  
 298 are going to be extracted (pressure, temperature, Mach, Reynolds and starting time of the  
 299 simulation). Therefore, the "Processing\_Step" class will be repeated once ("\_\_count\_\_": 1)  
 300 and will contain one "Name" and five "Parameter" and "Has\_numerical\_value" properties.  
 301 Regarding the "Tool" class, assume to separate between "Hardware" and "Software". Hence,  
 302 two instances of such a class ("\_\_count\_\_": 2) need to be created, each of them with its own  
 303 properties "System\_component", "Name" and "ID". The result of this manual file manipulation  
 304 is shown below.

```
305 {
306     "Processing_Step": {
```

```

3073         "__count__": 1,
3084         "__restrictions__": "",
3095         "Name": [1],
3106         "Parameter": [5],
3117         "Has_numerical_value": [5],
3123     },
3130     "Tool": {
3140         "__count__": 2,
3151         "__restrictions__": "",
3162         "System_component": [1,1],
3173         "Name": [1,1],
3184         "ID": [1,1],
3195     },
3206     ...
3217 }

```

**Listing 2: Filled-in .json file after step 2.**

322 Running `Multiplexer.py` expands the classes according to the parameters indicated in the  
323 previous step. The output is shown below. The new empty dictionary contains all the instances  
324 and corresponding properties the crawler will use in the creation of the metadata file.

```

3251 {
3262     "Processing_Step": {
3273         "__restrictions__": "",
3284         "Name": {
3295             "path": "",
3306             "type": "",
3317             "pattern": "",
3328             "postprocessor": {
3339                 "type": "",
3340                 "args": ""
3351             }
3362         },
3373         "Parameter_1": {
3384             "path": "",
3395             "type": "",
3406             "pattern": "",
3417             "postprocessor": {
3428                 "type": "",
3439                 "args": ""
3440             }
3451         }
3462     },
3473     ...

```

```

3484     "Has_numerical_value_1": {
3495         "path": "",
3506         "type": "",
3517         "pattern": "",
3528         "postprocessor": {
3539             "type": "",
3540             "args": ""
3551         }
3562     },
3573     ...
3584 },
3595 "Tool_1": {
3606     "__restrictions__": "",
3617     "Software_component": {
3628         "path": "",
3639         "type": "",
3640         "pattern": "",
3651         "postprocessor": {
3662             "type": "",
3673             "args": ""
3684         }
3695     },
3706     ...
3717 },
3728 ...
3739 }

```

**Listing 3:** Dictionary with all the classes and their properties expanded by the multiplexer in step 3.

374 The multiplexed dictionary has to be filled in manually again by the user. How to fill in the  
375 dictionary depends on how the user wants to retrieve the data and where the information is stored.  
376 In this example, data are all extracted from plain text files and the crawler uses regular expressions  
377 to locate and read the data. This is done by specifying the keywords: "path", "type" and  
378 "pattern". The entries in "postprocessor" can be left empty for the sake of this example.  
379 The lines below show how to hardcode metadata by providing a string in "type" (for the  
380 property "Parameter\_1", where the user directly provides the name of the parameter), retrieve  
381 information from a file using regular expressions ("Has\_numerical\_value\_1", retrieved from  
382 the file specified in "path") and from the output of a terminal command ("ID", where the  
383 terminal command is given in "pattern").

```

384 {
385     "Processing_Step": {
386         ...
387         "Parameter_1": {
388             "path": "",

```

```

389     "type": "string",
390     "pattern": "Freestream Mach number",
391     "postprocessor": {
392         "type": "",
393         "args": ""
394     }
395 }
396 },
397 ...
398 "Has_numerical_value_1": {
399     "path": "input.dat",
400     "type": "regex",
401     "pattern": "Mach :\\s(.*)\\n",
402     "postprocessor": {
403         "type": "",
404         "args": ""
405     }
406 },
407 ...
408 },
409 "Tool_1": {
410     ...
411     "ID": {
412         "path": "",
413         "type": "os",
414         "pattern": "hostname",
415         "postprocessor": {
416             "type": "",
417             "args": ""
418         }
419     },
420     ...
421 },
422 ...
423 }

```

**Listing 4:** Filled-in dictionary in step 4.

424 Finally, `EntityUtils.py` is used to run the actual extraction routine, which retrieves the meta-  
425 data according to the parameters specified in the previous step. The output file is shown below  
426 and could be either a `.json` or a `.yaml` file, according to the user needs.

```

427 {
428     "Processing_Step": {
429         "Name": "Wing simulation",

```

```
430         "Parameter_1": "Freestream Mach number",
431         "Parameter_2": "Freestream pressure",
432         "Parameter_3": "Freestream temperature",
433         "Parameter_4": "Freestream unit Reynolds number",
434         "Parameter_5": "Start time",
435         "Has_numerical_value_1": "0.35",
436         "Has_numerical_value_2": "61640",
437         "Has_numerical_value_3": "262",
438         "Has_numerical_value_4": "5.607E6",
439         "Has_numerical_value_5": "10:13:31"
440     },
441     "Tool_1": {
442         "System_component": "Hardware",
443         "Name": "lrz-coolmuc2-linux-cluster-2022"
444         "ID": "i22r07c05s05"
445     },
446     "Tool_2": {
447         "System_component": "Software",
448         "Name": "NSMB",
449         "ID": " 6.09.21    Date: 28 - January - 2021    "
450     },
451     "Method": {
452         "Name": "LU-SGS"
453     }
454 }
```

**Listing 5:** Final .json file containing the extracted metadata after step 5.

455 At this point, the generated metadata file can be stored together with the output data from the  
456 simulation. Whenever the user performs a new simulation and wants to extract the same type of  
457 metadata, there is no need to repeat all five steps of the process. In fact, the filled multiplexed  
458 dictionary created in step 4 will not change and acts as a configuration file that can be directly  
459 re-used in step 5. This means that the user needs only to add the command that runs step 5 in the  
460 script-based workflow. This corresponds to using HOMER in edge mode, which means invoking  
461 the crawler each time new data is generated at the end of a CFD simulation.

462 The other option would be to use the crawler after all the simulations have been run in order to  
463 retrieve all the metadata at once, which corresponds to using the crawler in central mode. In this  
464 case, the user will need to specify absolute file paths (via the keyword "path" in the multiplexed  
465 .json file) to point to the files containing the information. This means that the user needs to  
466 create an extra script that allows the crawler to search all the relevant folders and files. Such a  
467 script would be specialized according to the user's simulation environment and workflow. Hence,  
468 no example of such a usage can be given in the context of the generic CFD-showcase described  
469 in this work. However, an example script is provided in the GitLab folder for the Pizza-ontology  
470 tutorial. It must be noted that central and edge modes are not features of the tool itself, but are

471 different ways of using HOMER. It's up to the users to decide which is the best usage based  
472 on their needs and preferences. This, however, shows again the flexibility of the tool. **Another**  
473 **remark is that the user doesn't need to adhere strictly to the chosen ontology file, nor does the**  
474 **user have to use an ontology based on Metada4Ing. At any of the steps where manual input**  
475 **is needed, the user can adjust the classes and properties according to the case-specific needs.**  
476 **For example, the user could rename the property "Parameter" to "Variable" in the class**  
477 **"Processing\_Step" by manually amending the ".json" file while filling out the flat-class**  
478 **template during step 2. In fact, an ontology file is not even necessary for the actual extraction of**  
479 **the metadata, in principle. The user could even create it's own .json file with its own classes**  
480 **and properties skipping the first two steps altogether.** On one hand, of course, this approach  
481 requires a certain amount of overhead from the researcher side in terms of planning and preparing  
482 the .json files. On the other hand, it gives much more freedom to the user when it comes to  
483 adapt the crawler to the specific case at hand.

484 **Regarding the limitations of HOMER, the tool works best within standardized workflows, where**  
485 **the structure of the files containing the metadata to be extracted changes very little or not at**  
486 **all over time. Although, as shown, it would be possible to adapt the crawler, and in particular**  
487 **the multiplexed dictionary, to new file structures thanks to the flexibility of the tool, such an**  
488 **operation could take a considerable amount of time and effort from the user side if performed**  
489 **for every new application of a (changing) work flow. Hence, it appears sensible to limit the use**  
490 **of the crawler to cases where well-known and relatively fixed data structures are employed as it**  
491 **is common in most numerical and experimental research projects. The second limitation is the**  
492 **range of data formats the crawler can currently extract metadata from, which is limited to text**  
493 **and HDF5 files, together with outputs of terminal commands and hardcoded lines. Although the**  
494 **regular-expression parser allows to retrieve information from virtually any text file regardless**  
495 **of its extension, commonly used formats such as .xml have not been implemented, yet. As the**  
496 **crawler is designed flexible, this would be a straight forward process.**

## 497 5 Conclusion and Future Developments

498 In this work, HOMER (**H**PMC tool for **O**ntology-based **M**etadata **E**xtraction and **R**e-use), a tool  
499 to automate metadata extraction in script-based workflows, has been presented. The crawler,  
500 a python-written code, allows for a flexible approach to metadata retrieval. As starting point,  
501 the user can provide an ontology file, whose metadata scheme represents the backbone of the  
502 extracted information. The classes and attributes from the ontology can be tailored to the specific  
503 case at hand and expanded by means of the multiplexer. Once the user has filled in the final  
504 dictionary, the actual metadata extraction is executed. This can happen both in edge mode  
505 (natural application for script-based workflows) or, with some further user input, in central mode.  
506 Then, the extracted metadata can be further post-processed by some routines included in the  
507 code. The use of the tool requires some user input and tuning for the first application, but after  
508 that, it can be seamlessly integrated in potentially any workflow.

509 Currently, metadata can be retrieved from text and HDF5 files, from outputs of console commands  
510 or can be directly hardcoded in the configuration file. This limitation can be easily overcome  
511 in the future, as the **code is designed** in a modular way, thus allowing for a simple integration of

512 new building blocks. According to the user's needs, new readers/writers of other file formats  
513 can be added. The same applies for the post-processing capabilities on the extracted metadata.  
514 Moreover, work to increase the amount of readable file formats is planned, at first focusing on  
515 the most common formats in CFD applications.

516 As of now, HOMER can be already implemented in HPMC workflows, so as to enrich each  
517 processing step (e.g. mesh generation, simulation, post-processing, report) by adding the cor-  
518 responding metadata. This capability allows for the collection of valuable data (such as the  
519 energy consumption for a set of simulations) to enable secondary research and the development  
520 of new methodologies in HPC systems. In the current state, the tool would provide the best  
521 performance when used to extract metadata right after the creation of the data, as no parallel  
522 implementation for file parsing is present, yet. In its five-steps implementation described in this  
523 work, HOMER was mainly employed in the data life cycle for the processing stages of planning,  
524 creating/collecting and processing/analyzing. A future development of the tool would be to cover  
525 the complete data life cycle in a holistic approach, by providing the possibility to automatically  
526 preserve and publish/share the extracted metadata along with the research dataset. Through  
527 publishing not just the data but also administrative (preservation) metadata, third party users will  
528 be able to retrieve crucial information about accessibility, access rights and licenses among others.  
529 Bibliographic (author, identifiers) and descriptive (research domain, tools, methods, processing  
530 steps) metadata can be published in repositories together with the referenced research data, or  
531 be linked to the research data by persistent links and identifiers, if technical or organizational  
532 reasons impede a joint provisioning (for example, if the research data are too large to be stored  
533 in a common repository).

534 One main factor of making data FAIR is the use of a controlled vocabulary with common  
535 terminology. This is guaranteed by the fact that HOMER supports the usage of semantic  
536 ontologies as metadata schemes. These schemes have to be matched somehow with searchable  
537 metadata fields in the corresponding repositories. Only few repositories offer such publishing  
538 options, like DaRUS (University of Stuttgart) [24], which uses predefined metadata blocks, or  
539 Coscine (RWTH) [25], which provides the possibility to use standardized or self-created metadata  
540 application profiles [26]. These schemes still have to be parsed with the corresponding metadata  
541 fields in the extracted metadata file, to provide the metadata in a standardized, searchable and  
542 indexable front end. The NFDI4Ing consortium is simultaneously working on a generic interface  
543 which combines different kinds of metadata and data repositories with one standard-based  
544 interface. This enables the linking between all data and metadata of the research data life cycle,  
545 including experiments, raw data, software, subject-specific metadata sets, and the tracking of  
546 usage and citations. Standardized and automatically extracted metadata files can easily be  
547 made findable and accessible by this new generic interface [27]. Therefore, HOMER can be a  
548 crucial piece within the metadata toolchain from using common vocabularies and automatized  
549 extracting to FAIR publishing. The already mentioned Metadata4Ing ontology has been used  
550 as the reference during the early stages of the development of HOMER. In the meantime, a  
551 HPMC-sub-ontology has been developed within Metadat4Ing. Hence, one of the next steps will  
552 be to further adapt HOMER to this new sub-ontology, allowing the tool to be more effective in  
553 the complete data life cycle of a CFD workflow on HPC systems.

## 554 6 Acknowledgements

555 The authors would like to thank the Federal Government and the Heads of Government of the  
 556 Länder, as well as the Joint Science Conference (GWK), for their funding and support within  
 557 the framework of the NFDI4Ing consortium. Funded by the German Research Foundation  
 558 (DFG) - project number 442146713. Moreover, the authors gratefully acknowledge the Gauss  
 559 Centre for Supercomputing e.V. ([www.gauss-centre.eu](http://www.gauss-centre.eu)) for funding this project by providing  
 560 computing time on the GCS Supercomputer SuperMUC-NG at Leibniz Supercomputing Centre  
 561 ([www.lrz.de](http://www.lrz.de)).

## 562 7 Roles and contributions

563 **Giuseppe Chiapparino:** Conceptualization; Investigation; Methodology; Software - testing;  
 564 Validation; Writing – original draft

565 **Benjamin Farnbacher:** Data curation; Investigation; Writing – original draft (Introduction and  
 566 Conclusions)

567 **Nils Hoppe:** Conceptualization; Investigation; Methodology; Software - development and  
 568 design

569 **Radoslav Ralev:** Software - design, development, implementation and testing

570 **Vasiliki Sdralia:** Writing – original draft (Introduction)

571 **Christian Stemmer:** Funding acquisition; Resources; Supervision; Writing – review and editing  
 572 of original

## 573 References

- 574 [1] P. B. Heidorn, “Shedding Light on the Dark Data in the Long Tail of Science,” *Library*  
 575 *Trends*, vol. 57, no. 2, pp. 280–299, 2008. DOI: [doi:10.1353/lib.0.0036](https://doi.org/10.1353/lib.0.0036).
- 576 [2] B. Schembera and J. M. Duràn, “Dark Data as the New Challenge for Big Data Science  
 577 and the Introduction of the Scientific Data Officer,” *Philosophy & Technology*, vol. 33,  
 578 pp. 93–115, 2020. DOI: <https://doi.org/10.1007/s13347-019-00346-x>.
- 579 [3] NFDI4Ing Consortium. “Website.” (2022), [Online]. Available: <https://nfdi4ing.de>.
- 580 [4] Metadata4Ing Workgroup. “Metadata4ing: An ontology for describing the generation of  
 581 research data within a scientific activity.” (2022), [Online]. Available: [https://nfdi4i  
 582 ng.pages.rwth-aachen.de/metadata4ing/metadata4ing/index.html#ref](https://nfdi4ing.pages.rwth-aachen.de/metadata4ing/metadata4ing/index.html#ref).
- 583 [5] DCMI Usage Board. “Dcmi metadata terms.” (2020), [Online]. Available: [https://www  
 584 .dublincore.org/specifications/dublin-core/dcmi-terms/](https://www.dublincore.org/specifications/dublin-core/dcmi-terms/).
- 585 [6] Lebo, Timothy and Satya, Sahoo and Deborah, McGuinness. “Prov-o: The prov ontology.”  
 586 (2013), [Online]. Available: <https://www.w3.org/TR/prov-o/>.



- 587 [7] S. Liang, V. Holmes, G. Antoniou, and J. Higgins, "Icurate: A research data management  
588 system," in *Multi-disciplinary Trends in Artificial Intelligence*, A. Bikakis and X. Zheng,  
589 Eds., Cham: Springer International Publishing, 2015, pp. 39–47, ISBN: 978-3-319-26181-  
590 2. DOI: [10.1007/978-3-319-26181-2\\_4](https://doi.org/10.1007/978-3-319-26181-2_4).
- 591 [8] C. S. Adorf, P. M. Dodd, V. Ramasubramani, and S. C. Glotzer, "Simple data and workflow  
592 management with the signac framework," *Computational Materials Science*, vol. 146,  
593 pp. 220–229, 2018, ISSN: 0927-0256. DOI: <https://doi.org/10.1016/j.commatsci.2018.01.035>.
- 594 [9] L. Nagel and D. Lycklama, "Design principles for data spaces - position paper," version 1.0,  
595 2021. DOI: [10.5281/zenodo.5105744](https://doi.org/10.5281/zenodo.5105744).
- 596 [10] T. J. Skluzacek, "Dredging a data lake: Decentralized metadata extraction," in *Proceedings  
597 of the 20th International Middleware Conference Doctoral Symposium*, ser. Middleware  
598 '19, Davis, California: Association for Computing Machinery, 2019, pp. 51–53, ISBN:  
599 9781450370394. DOI: <https://doi.org/10.1145/3366624.3368170>.
- 600 [11] T. J. Skluzacek, R. Chard, R. Wong, *et al.*, "Serverless workflows for indexing large  
601 scientific data," in *Proceedings of the 5th International Workshop on Serverless Computing*,  
602 ser. WOSC '19, Davis, CA, USA: Association for Computing Machinery, 2019, pp. 43–48,  
603 ISBN: 9781450370387. DOI: <https://doi.org/10.1145/3366623.3368140>.
- 604 [12] J. Dixon, "Pentaho, hadoop, and data lakes." (2010), [Online]. Available: [https://jame  
605 sdixon.wordpress.com/2010/10/14/pentaho-hadoop-and-data-lakes/](https://jamesdixon.wordpress.com/2010/10/14/pentaho-hadoop-and-data-lakes/).
- 606 [13] D. Merkel, "Docker: Lightweight linux containers for consistent development and deploy-  
607 ment," *Linux Journal*, vol. 2014, no. 239, p. 2, 2014.
- 608 [14] I. Foster, "Globus online: Accelerating and democratizing science through cloud-based  
609 services," *IEEE Internet Computing*, vol. 15, no. 3, pp. 70–73, 2011. DOI: [10.1109  
610 /MIC.2011.64](https://doi.org/10.1109/MIC.2011.64).
- 611 [15] B. Allen, J. Bresnahan, L. Childers, *et al.*, "Software as a service for data scientists," *IEEE  
612 Internet Computing*, vol. 55, no. 2, pp. 81–88, 2012. DOI: [10.1145/2076450.2076468](https://doi.org/10.1145/2076450.2076468).
- 613 [16] B. Schembera, "Like a rainbow in the dark: Metadata annotation for HPC applications in  
614 the age of dark data," *Journal of Supercomputing*, vol. 77, pp. 8946–8966, 2021. DOI:  
615 <https://doi.org/10.1007/s11227-020-03602-6>.
- 616 [17] B. Schembera and D. Iglezakis, "The Genesis of EngMeta - A Metadata Model for  
617 Research Data in Computational Engineering," in *Metadata and Semantic Research*,  
618 Cham: Springer International Publishing, 2019, pp. 127–132, ISBN: 978-3-030-14401-2.  
619 DOI: [https://doi.org/10.1007/978-3-030-14401-2\\_12](https://doi.org/10.1007/978-3-030-14401-2_12).
- 620 [18] S. Padhy, G. Jansen, J. Alameda, *et al.*, "Brown dog: Leveraging everything towards  
621 autocuration," in *2015 IEEE International Conference on Big Data (Big Data)*, 2015,  
622 pp. 493–500. DOI: [10.1109/BigData.2015.7363791](https://doi.org/10.1109/BigData.2015.7363791).
- 623 [19] G. P. Rodrigo, M. Henderson, G. H. Weber, C. Ophus, K. Antypas, and L. Ramakrishnan,  
624 "ScienceSearch: Enabling search through automatic metadata generation," in *2018 IEEE  
625 14th International Conference on e-Science (e-Science)*, 2018, pp. 93–104. DOI: [10.110  
626 9/eScience.2018.00025](https://doi.org/10.1109/eScience.2018.00025).
- 627

- 628 [20] K. Frey, K. Schneider, O. Maus, and T. Mühlhaus. “Swate: A swate workflow annotation  
629 tool for excel.” (2022), [Online]. Available: [https://github.com/nfdi4plants/Swa](https://github.com/nfdi4plants/Swate)  
630 [te](https://github.com/nfdi4plants/Swate).
- 631 [21] UK Data Service, modified by TUM University Library (UB). “Data life cycle - icons.”  
632 (2022).
- 633 [22] M. A. Musen, “The protégé project: A look back and a look forward,” *AI Matters*, vol. 1,  
634 no. 4, pp. 4–12, 2015. DOI: [10.1145/2757001.2757003](https://doi.org/10.1145/2757001.2757003). [Online]. Available: <https://doi.org/10.1145/2757001.2757003>.
- 636 [23] J. Vos, N. Duquesne, and H. J. Lee, “Shock wave boundary layer interaction studies using  
637 the NSMB flow solver,” in *3rd European Symposium on Aerothermodynamics for Space*  
638 *Vehicles*, ESA SP-426, 1999.
- 639 [24] University of Stuttgart. “Darus.” (2022), [Online]. Available: [https://www.izus.uni-](https://www.izus.uni-stuttgart.de/en/fokus/darus/)  
640 [stuttgart.de/en/fokus/darus/](https://www.izus.uni-stuttgart.de/en/fokus/darus/).
- 641 [25] RWTH Aachen University. “Coscine.” (2022), [Online]. Available: [https://coscine](https://coscine.rwth-aachen.de)  
642 [.rwth-aachen.de](https://coscine.rwth-aachen.de).
- 643 [26] RWTH Aachen University. “Aims – applying interoperable metadata standards.” (), [On-  
644 line]. Available: [https://www.wzl.rwth-aachen.de/cms/wzl/Forschung/Forsch](https://www.wzl.rwth-aachen.de/cms/wzl/Forschung/Forschungsumfeld/Forschungsprojekte/Projekte/~ivong/ProMiDigit-Process-Mining-fuer-No-Code/)  
645 [ungsumfeld/Forschungsprojekte/Projekte/~ivong/ProMiDigit-Process-Mi](https://www.wzl.rwth-aachen.de/cms/wzl/Forschung/Forschungsumfeld/Forschungsprojekte/Projekte/~ivong/ProMiDigit-Process-Mining-fuer-No-Code/)  
646 [ning-fuer-No-Code/](https://www.wzl.rwth-aachen.de/cms/wzl/Forschung/Forschungsumfeld/Forschungsprojekte/Projekte/~ivong/ProMiDigit-Process-Mining-fuer-No-Code/).
- 647 [27] NFDI4Ing Consortium. “Metadata hub.” (2022), [Online]. Available: [https://git.rwth](https://git.rwth-aachen.de/nfdi4ing/s-3/s-3-3/metadatahub)  
648 [aachen.de/nfdi4ing/s-3/s-3-3/metadatahub](https://git.rwth-aachen.de/nfdi4ing/s-3/s-3-3/metadatahub).