

# Guidance framework to apply best practices in ecological data analysis: Lessons learned from building Galaxy-Ecology

Royaux Coline<sup>1,2\*</sup>, Mihoub Jean-Baptiste<sup>3</sup>, Jossé Marie<sup>4</sup>, Pelletier Dominique<sup>5</sup>, Norvez Olivier<sup>6</sup>, Reecht Yves<sup>7,8</sup>, Fouilloux Anne<sup>9</sup>, Rasche Helena<sup>10</sup>, Hiltemann Saskia<sup>11</sup>, Batut Bérénice<sup>12,13</sup>, Eléaume Marc<sup>14,15</sup>, Segumineau Pauline<sup>14,15</sup>, Massé Guillaume<sup>16</sup>, Amossé Alan<sup>17</sup>, Bissery Claire<sup>8,18</sup>, Lorrilliere Romain<sup>3</sup>, Martin Alexis<sup>19</sup>, Bas Yves<sup>3,20</sup>, Virgoulay Thimothée<sup>21,22</sup>, Chambon Valentin<sup>17</sup>, Arnaud Elie<sup>2</sup>, Michon Elisa<sup>23</sup>, Urfer Clara<sup>2,24</sup>, Trigodet Eloïse<sup>21,24</sup>, Delannoy Marie<sup>3</sup>, Loïs Gregoire<sup>3</sup>, Julliard Romain<sup>3</sup>, Grüning Björn<sup>25</sup>, The Galaxy-E community, Le Bras Yvan<sup>2</sup>

<sup>1</sup> UMR8067 Biologie des Organismes et Ecosystèmes Aquatiques (BOREA, MNHN-CNRS-SU-IRD-UCN-UA), Sorbonne Université, Station Marine de Concarneau - Concarneau, France

<sup>2</sup> Pôle national de données de biodiversité, UAR2006 PatriNat (OFB-MNHN-CNRS-IRD), Muséum National d'Histoire Naturelle, Station Marine de Concarneau - Concarneau, France

<sup>3</sup> Centre d'Écologie et des Sciences de la Conservation (UMR7204 CESCO, MNHN-CNRS-SU), Muséum National d'Histoire Naturelle, Sorbonne Université, Centre National de la Recherche Scientifique - Paris, France

<sup>4</sup> Data Terra, Centre National de la Recherche Scientifique - Brest, France

<sup>5</sup> UMR DECOD (Ifremer-Agrocampus Ouest-INRAE) - Lorient, France

<sup>6</sup> Pôle National de Données de Biodiversité, UAR2006 PatriNat (OFB-MNHN-CNRS-IRD), Fondation pour la Recherche sur la Biodiversité, Muséum national d'Histoire naturelle - Paris, France

<sup>7</sup> Institute of Marine Research - Bergen, Norway

<sup>8</sup> Institut français de recherche pour l'exploitation de la mer (Ifremer) - Brest, France

<sup>9</sup> Simula Research Laboratory - Oslo, Norway

<sup>10</sup> Department of Pathology and Clinical Bioinformatics, Erasmus Medical Center - Rotterdam, The Netherlands

<sup>11</sup> Institute of Pharmaceutical Sciences, Faculty of Chemistry and Pharmacy, University of Freiburg - Freiburg, Germany

<sup>12</sup> Institut Français de Bioinformatique, CNRS UAR3601 - Évry, France

<sup>13</sup> Mésocentre, Clermont-Auvergne, Université Clermont Auvergne - Clermont-Ferrand, France

<sup>14</sup> Institut de Systématique Evolution, Biodiversité (UMR7205 ISYEB, MNHN-CNRS-SU-EPHE), Département Origines et Évolution, Muséum national d'Histoire naturelle - Paris, France

<sup>15</sup> Institut de Systématique Evolution, Biodiversité (UMR7205 ISYEB, MNHN-CNRS-SU-EPHE), Département Origines et Évolution, Station Marine de Concarneau - Concarneau, France

<sup>16</sup> UMR LOCEAN (CNRS-SU-IRD-MNHN), Centre National de la Recherche Scientifique, Station Marine de Concarneau - Concarneau, France

<sup>17</sup> Muséum National d'Histoire Naturelle, Station Marine de Concarneau - Concarneau, France

<sup>18</sup> Université Claude Bernard Lyon 1 - Lyon, France

43 <sup>19</sup> UMR8067 Biologie des Organismes et Ecosystèmes Aquatiques (BOREA, MNHN-CNRS-SU-IRD-UCN-UA),  
44 Muséum national d'Histoire naturelle - Paris, France

45 <sup>20</sup> UAR2006 PatriNat (OFB-MNHN-CNRS-IRD), Muséum national d'Histoire naturelle - Paris, France

46 <sup>21</sup> Centre d'Écologie et des Sciences de la Conservation (UMR7204 CESCO, MNHN-CNRS-SU), Muséum  
47 National d'Histoire Naturelle - Concarneau, France

48 <sup>22</sup> Université de Montpellier - Montpellier, France

49 <sup>23</sup> Institut des Sciences de la Mer de Rimouski, Université du Québec à Rimouski - Rimouski, Québec, Canada

50 <sup>24</sup> Université de Bretagne Occidentale - Brest, France

51 <sup>25</sup> Bioinformatics Group, Department of Computer Science, Albert-Ludwigs-University Freiburg -  
52 Freiburg, Germany

53

54 \*Corresponding author

55 Correspondence: coline.royaux@mnhn.fr

56

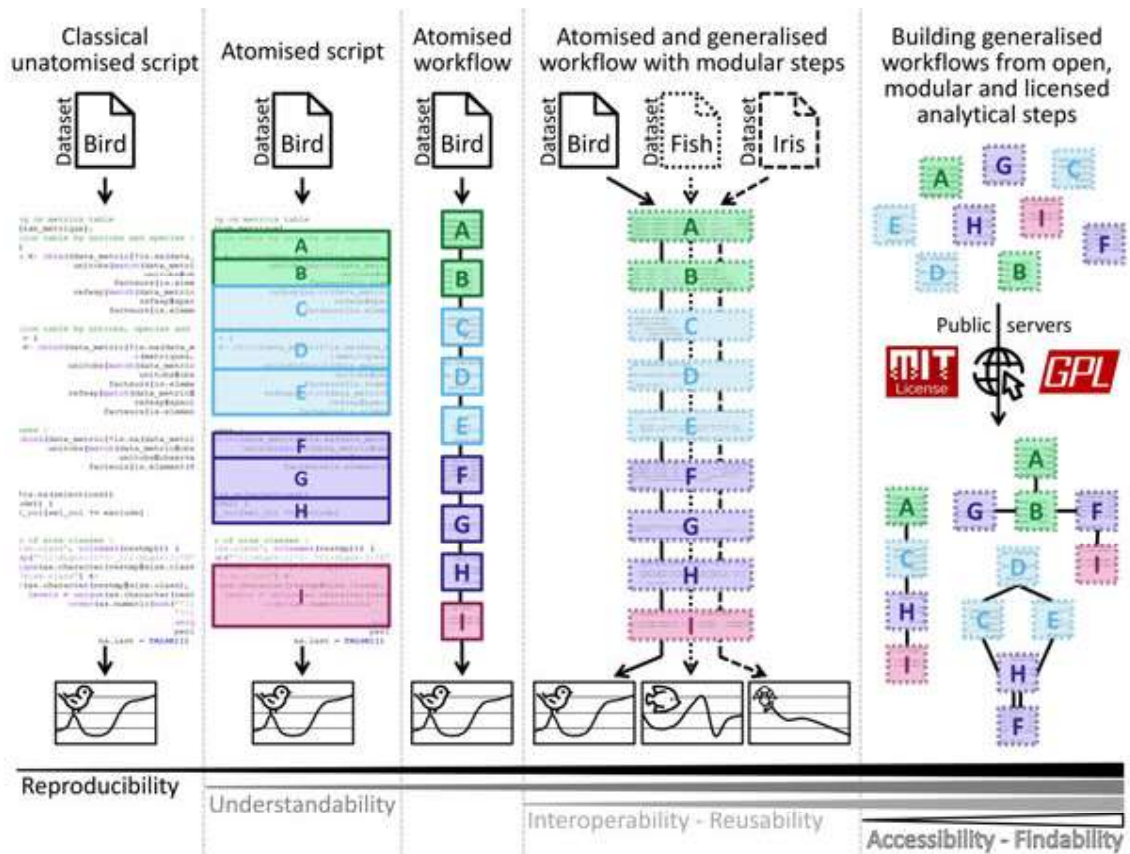
## 57 **ABSTRACT**

58 Numerous conceptual frameworks exist for best practices in research data and analysis  
59 (*e.g.* Open Science and FAIR principles). In practice, there is a need for further progress to  
60 improve transparency, reproducibility, and confidence in ecology. Here, we propose a  
61 practical and operational framework for researchers and experts in ecology to achieve  
62 best practices for building analytical procedures from individual research projects to  
63 production-level analytical pipelines. We introduce the concept of atomisation to identify  
64 analytical steps which support generalisation by allowing us to go beyond single analyses.  
65 The term atomisation is employed to convey the idea of single analytical steps as “atoms”  
66 composing an analytical procedure. When generalised, “atoms” can be used in more than  
67 a single case analysis. These guidelines were established during the development of the  
68 Galaxy-Ecology initiative, a web platform dedicated to data analysis in ecology. Galaxy-  
69 Ecology allows us to demonstrate a way to reach higher levels of reproducibility in  
70 ecological sciences by increasing the accessibility and reusability of analytical workflows  
71 once atomised and generalised.

72

73  
74

**Graphical abstract** – Levels of attainable best practices through the atomisation  
– generalisation framework



75

76 **Keywords:** Biodiversity; Reproducible analyses; Galaxy; Best practices; Atomisation;  
77 Generalisation; Workflows; Ecoinformatics; Conda; Container; Common Workflow Language; RO-  
78 CRATE

## 80 Ecology's Reproducibility Crisis

81 Research in ecology is increasingly shaped by the availability of novel analytical solutions  
82 and statistical tools. Given the ever-growing amount of data available, much attention is often  
83 given to the thought process behind statistical analyses to handle different data distributions,  
84 pseudo-replication, and sampling biases for instance (*NERC* 2010, 2012; Hampton *et al.*, 2017;  
85 Emery *et al.*, 2021). Despite the high-quality standards required by the scientific community  
86 from data access to analysis, the level of complexity of ecological systems makes results  
87 difficult to reproduce. The ongoing "reproducibility crisis" has also led researchers to pay  
88 closer attention to the quality of analyses to increase confidence in their studies and  
89 conclusions (Ioannidis, 2022; Fanelli, 2018). Reproducibility (*i.e.* different teams and  
90 experimental setups obtaining similar results; Plesser, 2018) is one of the main criteria for  
91 evaluating robust science and reliable conclusions. The term "reproducibility" is a relative  
92 concept and has known various definitions depending on field and context. Reproducibility of  
93 analyses ("computational reproducibility") is defined by Cohen-Boulakia *et al.* (2017) as the  
94 ability of distinct analyses to reach to the same conclusion.

95 In the current context of the global biodiversity crisis, the scientific community needs to  
96 use all available data and provide as robust as possible evidence regarding the state and  
97 dynamic of ecological systems, from genetic to ecosystem. At the same time, using analytical  
98 tools to provide robust evidence can be complex and may require advanced skills that are not  
99 widely available across the scientific community (Hampton *et al.*, 2017). Therefore,  
100 operational solutions and methodological guidelines can allow analytical workflows to be  
101 more accessible without degrading the scientific quality of analyses, and thus, promote  
102 efficient and broad deployment of best practices.

## 103 Is the ecology community failing to meet best practices?

104 The first step towards reproducibility is knowing current best practices and  
105 recommendations. Among them, the FAIR principles (Wilkinson *et al.*, 2016), for which the  
106 availability of the data and the code used for each published result is an essential criterion,  
107 may be key for appropriate management through the data life cycle (Michener, 2015). The  
108 FAIR principles (see also CARE principles by Carroll *et al.*, 2020) are considered as a founding  
109 framework to share data along four important elements: "Findable" for humans and  
110 machines; "Accessible" with a detailed access procedure; "Interoperable" for interaction with  
111 other data or applications; "Reusable" in an identical or different context. In addition to these  
112 principles, propositions have been delimited within several thematic communities in ecology  
113 to evaluate and enhance best practices application, notably the Species Distribution  
114 Modelling communities (Araújo *et al.*, 2019; Zurell *et al.*, 2020).

115 Although data accessibility has been substantially improved in ecology during the past  
116 decade, sharing analytical scripts and codes remain largely marginal (Archmiller *et al.*, 2020;  
117 Culina *et al.*, 2020; Minocher *et al.*, 2021; Ivimey-Cook *et al.*, 2023). However, even if sharing  
118 code is necessary to achieve good computational reproducibility, it is insufficient. Therefore,  
119 the utilisation of computational workflows has been suggested as a solution for improving  
120 computational reproducibility (Cohen-Boulakia *et al.*, 2017; Grüning *et al.*, 2018) through

121 software such as Snakemake (Köster & Rahmann, 2012), Nextflow (Di Tommaso *et al.*, 2017),  
122 or Galaxy (The Galaxy Community, 2022). A workflow is generally defined as a sequence of  
123 distinct computational tasks for a particular objective (Goble *et al.*, 2020). As such, a workflow  
124 represents the backbone of a single specific analysis. Throughout the analytical procedure, a  
125 typical workflow starts with raw data, which can be extracted from several databases or data  
126 files and processed through a series of analytical steps. The products resulting from these  
127 analytical steps (*i.e.* the outputs of the computational workflow) can be data files, graphic  
128 representations and any associated metrics.

129 When properly designed, a certain level of reproducibility can be easily achieved since  
130 workflow languages naturally capture the following four key elements (Cohen-Boulakia *et al.*,  
131 2017):

- 132 – the specificities of the workflow, the analysis steps and associated tools;
- 133 – the workflow entries, datasets and parameters;
- 134 – the environment and context of the use of the workflow;
- 135 – the results obtained and the outputs of the workflow.

136 In the original publication of Wilkinson *et al.* (2016), the focus of FAIR principles was  
137 mainly on observational data. However, the principles can be applied to software and  
138 computational workflows (Lamprecht *et al.*, 2019; Goble *et al.*, 2020). For instance, a code  
139 shared as supplementary material of a non-open access publication could be considered as  
140 "Interoperable" but is not easily "Findable", "Accessible", or "Reusable". In contrast, a large  
141 block of code consisting of several hundred lines, from data pre-processing to final results and  
142 graphics as pictured in the Graphical abstract ①, may require efforts to understand and  
143 adapt to other kinds of data ("non-reusable"), mainly if annotations or comments are limited.  
144 Similarly, an analytical procedure shared without indicating the versions of hardware,  
145 software, and packages has a low chance of producing identical outputs, making it less  
146 reproducible. These issues may harm the scientific community by preventing fully transparent  
147 communication among users about knowledge production and practice comparison. They  
148 can also be detrimental to individual authors, when they need to update or run new analyses.

## 149 **Impact on Ecology Research**

150 The efficiency of the scientific process is greatly affected by the lack of computational  
151 reproducibility and FAIRness of analytical procedures. The adoption of FAIR practices was  
152 estimated to save 10.2 billion € per year in Europe (Munafò *et al.*, 2017; European commission,  
153 2018; Gomes *et al.*, 2022). Moreover, consistent application of reproducibility and FAIR  
154 principles will improve trust in research studies and scientific reports (Powers & Hampton,  
155 2019; Lortie, 2021; Jenkins *et al.*, 2023).

156 The widespread use of computational languages to process large-scale data and analyse  
157 complex systems has been a major advance in studying the ecosphere at any spatio-temporal  
158 scale (Michener & Jones, 2012; Farley *et al.*, 2018). However, the ever-growing technical and  
159 programming skills required to take advantage of such computational solutions by the  
160 scientific community raise new challenges (Jetz *et al.*, 2019; Leroy, 2022; Boyd *et al.*, 2023).  
161 The use of increasingly complex analytical solutions, paired with different approaches or  
162 programming languages, creates barriers to uptake and challenges for peer-review. Indeed, many  
163 ecologists have acquired their programming skills through self-study or through courses that

164 combine instruction in statistics and ecological principles with an introduction to  
165 programming. This learning process does not inherently compromise the quality of the  
166 analyses and results; however, it may lead to inappropriate coding habits. As a response to  
167 this situation, adequate training was identified by life science researchers (*Community Survey*  
168 *Report*, 2013; Williams & Teal, 2017; Larcombe *et al.*, 2017), as it would help involve more  
169 people in the understanding of current analytical solutions and benefit to scientific  
170 cooperation (Touchon & McCoy, 2016; Gownaris *et al.*, 2022). Research is typically structured  
171 through a highly competitive organisation, with a potentially detrimental effect on scientific  
172 knowledge (Fang & Casadevall, 2015). Instead, fostering collaboration and collective  
173 intelligence by promoting transparent sharing of analytical procedures, would offer more  
174 persistent and robust ways to achieve actionable science (Ellemers, 2021). Such efforts would  
175 be of paramount importance in environmental sciences and the conservation of biodiversity  
176 by providing governance and guiding actions with increasingly robust evidence (Keenan *et al.*,  
177 2012).

## 178 **Are there simple and ready-to-use solutions?**

179 In this note, we aim to promote the reuse of existing concepts and solutions as pillars  
180 toward better practices for ecological analyses by providing a streamlined framework. We  
181 believe the atomisation-generalisation framework presented in the second part of this note  
182 represents an operational and actionable path for researchers and experts to attain levels of  
183 best practices (*e.g.* reproducibility, FAIR, open science, R compendium; Casajus N., 2023) with  
184 no more investment than they are able or willing to provide (Field *et al.*, 2014). Atomisation is  
185 used to refer to the identification of distinct analytical steps each constituting an analytical  
186 procedure. It is a non-standard term introduced in this note to convey the idea of analytical  
187 “atoms”. As for atom particles that etymologically correspond to “indivisible” but are  
188 composed of subatomic particles, an analytical atom represents a single analytical step  
189 composed of several functions. Generalisation involves the alteration of an analytical step to  
190 enlarge its applicability in diverse contexts and for diverse purposes. Therefore, generalisation  
191 cannot be efficiently achieved without prior atomisation.

192 Atomisation and Generalisation are central organising principles in the design of the  
193 Galaxy-Ecology (Galaxy-E) initiative (see section III). Galaxy-E is a demonstration platform for  
194 applying best practices such as the FAIR principles and computational reproducibility for  
195 analytical procedures in ecology. Hence, this technical note is partly Galaxy-oriented, not to  
196 present the platform as a prescriptive solution but to give an operational example of the best  
197 practices it helps to achieve.

## 198 **Guidelines for best practices**

### 199 **Atomisation: what is it and why?**

200 Atomisation refers to dividing an analytical procedure into several specific steps (“atoms”;  
201 Graphical abstract ②) generating a suite of elementary analytical steps as pictured in the  
202 Graphical abstract ③. For instance, in a maximally-atomised workflow, each small step  
203 would be conducted by its own bespoke function. Breaking down the analytical process into  
204 atoms functioning as building blocks allows for better understanding, modularity, and

205 visibility of the analytical flow. It permits making it more accessible to a broader audience or  
206 facilitating the peer-review process. Indeed, an extended one-block code that imports raw  
207 data, makes pre-processing steps (*e.g.* filter, formatting), conducts analyses (*e.g.* distribution  
208 study, modelling), and performs final representations of results (*e.g.* maps, plots) can be  
209 challenging to understand and reuse by others or even the same person after some time.

210 McIntire *et al.* (2022) described the PERFICT approach (Prediction, Evaluation, Reusability,  
211 Free access, Interoperability, Continuous workflows, and routine Tests) to set a new  
212 foundation for models in predictive ecology. This can be applied more generally to the  
213 analytical procedure in ecology and biodiversity. In their article, McIntire and collaborators  
214 make an analogy between code development and Lego® construction, similar to our definition  
215 of atomisation. Functions are a workflow's most fundamental analytical steps and can be  
216 seen as modular pieces, alike single pieces of Lego®. Modules can be created from a single or  
217 series of successive functions comparably as in Lego® structures made of several pieces (*e.g.*  
218 meant to build cars, houses, or road). These modules (or atoms, tools) can be used as  
219 standalone or combined to make simple to complex analytical workflows (*e.g.* data  
220 formatting or curation, running statistical models, or generating graphical elements for  
221 visualisation). Doing so, the atomisation approach may facilitate sharing or teaching  
222 analytical practices since beginners can easily understand the general organisation of the  
223 analytical procedure by simply reading the list of steps in the analysis with a limited degree of  
224 complexity. Decoupling programming skills from analytical skills can make data processing  
225 more accessible to a wider audience. Indeed, once each elementary step is clearly identified  
226 and delimited along the atomisation process, it is easier to grasp the whole analytical  
227 procedure and focus on the review of each step at a time or (re)use it. New workflows can  
228 further be generated by recombining existing, validated or peer-reviewed elementary steps in  
229 innovative ways. This process can save time, increase confidence, and avoid potential  
230 programming mistakes, allowing greater focus on understanding the analytical workflow.

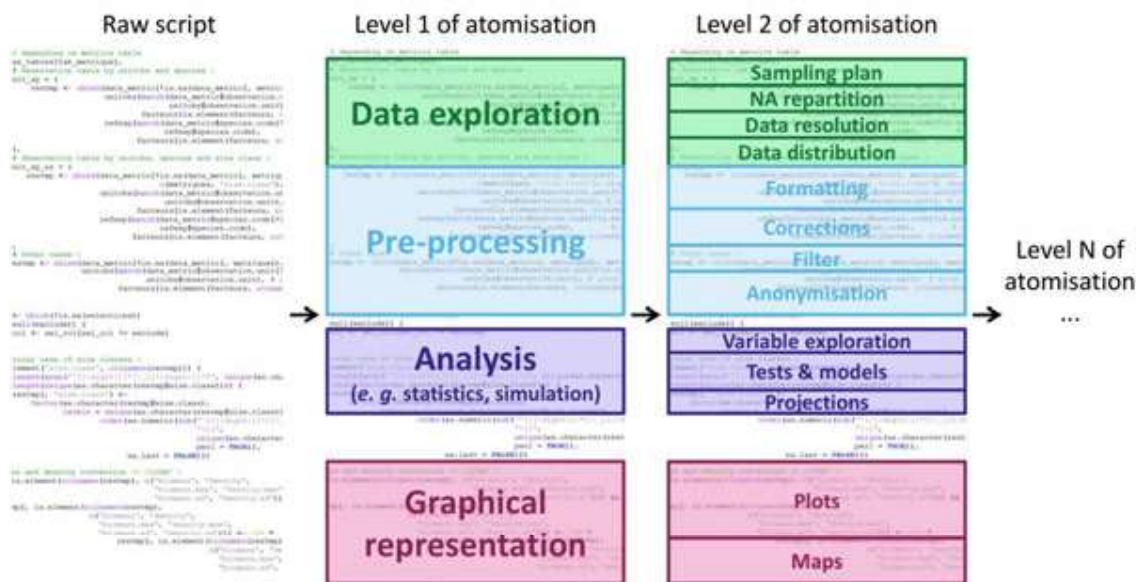
### 231 **Generalisation: what is it and why?**

232 Generalisation refers to the modification of an analytical procedure to make it applicable  
233 to many settings, by removing specificities related to a particular data file or data format. This  
234 means trying to avoid hard-coding anything that is specific to the structure of the original  
235 dataset (*e.g.* number of years). Generalisation aims to optimise the reusability at different  
236 times (*e.g.* regular result update), enlarge the application of a given analysis to different input  
237 data files while keeping the initial analytical procedure fully reproducible as pictured in the  
238 Graphical abstract ④. Generalising an analytical step requires identifying key elements and  
239 invariant parameters from those that must be adaptable to allow for the analysis to be  
240 applied to specific characteristics of various datasets. These parameters must be  
241 implemented to be easily modified if needed. Generalisation can be tricky because the higher  
242 the flexibility of an analytical step, the greater the risk of errors in its use. This is why  
243 generalisation should be complemented by clear statement and an implementation of red  
244 flags and warnings to prevent such events. As with atomisation, generalisation is primarily a  
245 conceptual way to build analytical procedures. It requires minor change of practices to reach  
246 certain degree of generalisation, avoiding additional effort later for reusability, reproducibility,  
247 and share.



248 **Practical steps towards atomised and generalised coding**

249 Breaking down codes into elementary steps to achieve atomisation is not an intuitive task  
250 at first as it may target a single function or a more intricate set of several functions. There  
251 could be different degrees of atomisation, depending on the grain required to decompose the  
252 analytical process (fig. 1; tab. 1). The application of general guidelines and best practices  
253 implies finding a balance between the most appropriate degree of atomisation and  
254 generalisation. This depends on the type of analytical procedure or the targeted audience (*e.g.*  
255 with different interests and programming skills). Attention to this balance is critical to ensure  
256 that the analytical procedures could be reused. For instance, a workflow in which each  
257 function would be considered as a unique elementary step would optimise the flexibility but  
258 may likely add unnecessary complexity. At the other extreme, considering a whole analytical  
259 workflow as an elementary step may make it ready-to-use and simplify its application, but  
260 would be too coarse and therefore limit flexibility by violating the principle of atomisation.



261

262 **Figure 1** - Illustration of the atomisation of an existing code. The first level of  
263 atomisation is delimitating the large sections of an analytical procedure that  
264 exist in almost all procedures. This first level is conveyed using same colours to  
265 the second level of atomisation where more detailed and specific analytical  
266 steps are illustrated in each section. The process of atomisation can continue  
267 through a multitude of levels, ultimately leading to the maximally atomised  
268 procedure, which is comprised of a single function.  
269



**Table 1** - Example of atomisation levels

Level 1 - big shape	Level 2	Level 3
Data exploration	Sampling plan	Complete Balanced
	Missing values	Proportion Distribution
	Data granularity	Geographic resolution Temporal resolution Measure resolution
	Data distribution	Geographic coverage Temporal coverage Measures ranges Summaries
...	...	...
Pre-processing	Formatting	Change file format Change general format
	Corrections	Remove special characters Remove low trust observations Correct measures
	Filtering	Remove unwanted observations
	Anonymisation	Anonymise names Anonymise localities Anonymise species
...	...	...
Analysis	Variable exploration	PCA Collinearity Correlation
	Unimodal tests	Linear Models $\chi^2$ Student
	Statistical models	Generalised Linear Models Generalised Additive Models Random Forest
	Models Evaluation	Evaluation metrics ( <i>e.g.</i> AIC, Jaccard) Validation methods
	Projections	Geographical projections Temporal projections
...	...	...
Representation	Plot	Raw variables Modelled results
	Map	Observations Projections
...	...	...

271 A few changes in code-writing habits can enhance the reusability of the analytical  
272 procedure by generating easy-to-understand analytical procedure without investing much  
273 time. It is best to develop each elementary step directly in separate code files and to give  
274 details of the order in which elementary steps are used for each analytical workflow. To  
275 ensure reproducibility and traceability of the results, each computation of the analytical  
276 workflow should be associated with the details of the parameters settings and datasets used.  
277 From a practical point of view, a couple of recommendations could be made for coding  
278 elementary steps to facilitate generalisation and ease the reuse. Once each elementary step is  
279 defined, we recommend all dependencies (*e.g.* software version, packages, libraries and their  
280 versions) to be set at the same place, at the start of the code, followed by modular parameters  
281 (*e.g.* input file location and name, column selection, modelling parameters, data specificities,  
282 output saving location). When the script of the elementary step is completed, modular  
283 parameters should be the only part of the code that may be modified in future reuse.  
284 Dependencies and subsequent computational tasks should be left untouched to ensure the  
285 integrity of the analysis and then, reproducibility. In the end, it is best to add an open-source

286 license to any analytical procedure shared publicly (e.g. MIT, GPL). It permits to clearly state  
287 the terms and conditions of diffusion, share and reuse.

288 As such, atomisation and generalisation may overcome social or psychological barriers  
289 related to transparent sharing, either related to securing ownership (e.g. DOI) and to  
290 embarrassment or fear during a peer-review process (Gomes *et al.*, 2022). Indeed, as  
291 atomisation and generalisation notably permit higher readability of codes, it would be more  
292 straightforward for the writer or even trusted peers to verify and review the steps before  
293 submission.

294 Atomisation and generalisation are related and complementary concepts that may be  
295 applied from the earliest stages of the programming development. Indeed, atomisation into  
296 adequate elementary steps is necessary to properly generalise an analytical procedure as it  
297 permits to enhance the modularity of the procedure and its capacity to be tailored to different  
298 data types.

### 299 **Entering a new dimension: the Galaxy-E initiative example**

300 Developing open and properly atomised and generalised analytical procedures can  
301 already represent a significant step forward in terms of best practice. Galaxy is a good  
302 illustration of atomisation and generalisation with easier management of analytical  
303 workflows. The platform proposes many analytical tools that represent generalised and  
304 atomised elementary steps. These tools are modular and openly licensed, which permits to  
305 build generalised workflows as pictured in the Graphical abstract 5.

306 Galaxy (The Galaxy Community, 2022) is a workflow-oriented web platform for analysing  
307 data and sharing outputs. It allows scientists to share, develop, and use various datasets and  
308 data processing tools (e.g. data formatting, statistical tests, graphic representations).

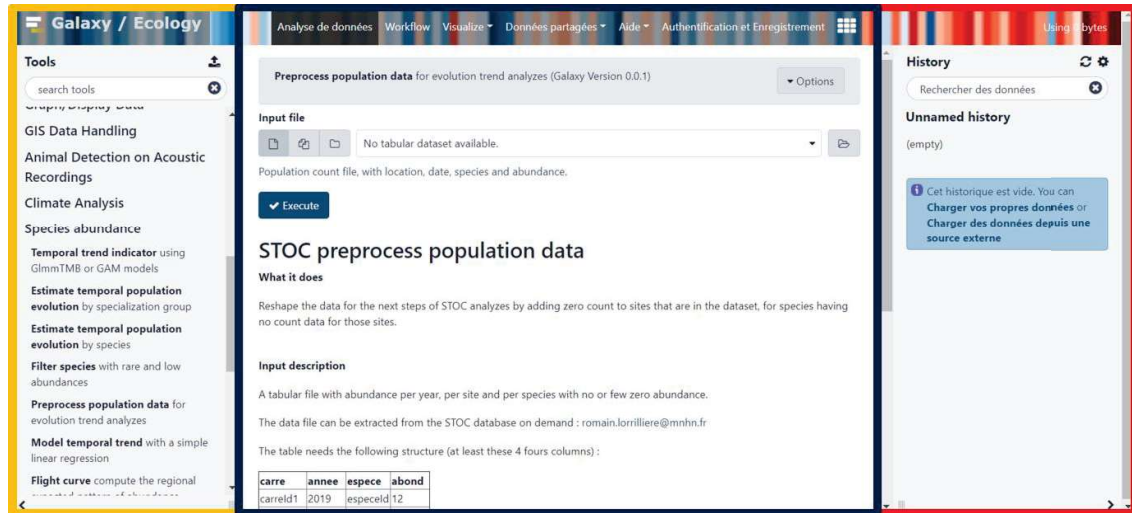
309 Galaxy enables good reproducibility for data exploration and analyses, helps compute  
310 intricate analyses on big data files, enables collaboration, and can support the teaching  
311 process. Galaxy-E is a Galaxy server dedicated to ecological analyses maintained by the  
312 European Galaxy team (supported by the German Federal Ministry of Education and Research  
313 and the German Network for Bioinformatics Infrastructure), and is available at  
314 <https://ecology.usegalaxy.eu>.

315 Galaxy-E is mostly aimed at scientists that process biodiversity data and already  
316 understand the general functioning of the analytical procedures they want to produce. The  
317 rationale for a user would be to create or reuse analytical workflows with high FAIRness in a  
318 collaborative and open source platform. It can be used for individual analyses as well as for  
319 collaborative projects. In some cases, if the analytical procedure is already clearly defined, it  
320 can be used by citizens or for teaching.

321 There are different Galaxy servers, at global, continental, and national levels (European  
322 and French levels for example), but also according to the fields (e.g., biomedical, ecology,  
323 climate). The Galaxy-E initiative is hosted by European (<https://ecology.usegalaxy.eu>) and  
324 French (<https://ecology.usegalaxy.fr>) servers.

325 Datasets can be uploaded on a Galaxy server from a local device, an online server, or a  
326 database. Users can then access every available tool (fig. 2, left panel) to modify, explore, and  
327 analyse their data. All tools used, parameters, and data (inputs and outputs) of the analysis  
328 are saved in a private “Galaxy history” (fig. 2, right panel), documenting every step of the

329 analytical procedure and recording the provenance of each output. From any history, the user  
330 can extract a workflow (fig. 3) or directly share or publish the history itself. Workflows are  
331 reusable through WorkflowHub (<https://workflowhub.eu>) or Dockstore (<https://dockstore.org>)  
332 and exportable in CWL and RO-CRATE standards.



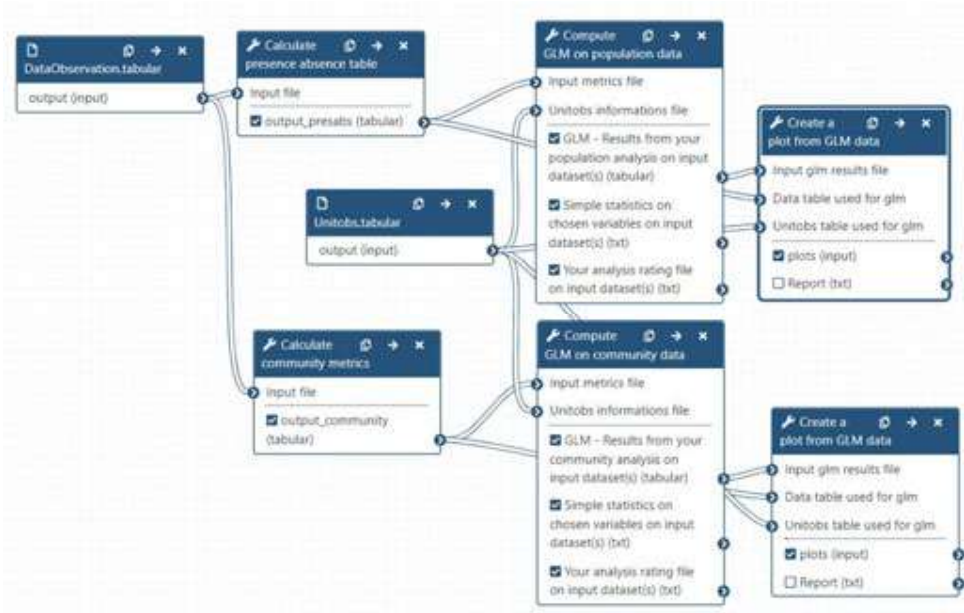
333

334

335

336

**Figure 2** - Galaxy-Ecology users' interface <https://ecology.usegalaxy.eu>. Yellow panel on the left: analysis tool list; blue panel in the middle: current tool interface; red panel on the right: Galaxy analysis history



337

338

339

340

341

342

**Figure 3** - Representation of a Galaxy workflow in the editing interface of a Galaxy server. Each box represents an analysis tool, and the lines represent the flow of data through the tools. In relation with the atomisation-generalisation framework, each box (tool) corresponds to an atomised and generalised step with editable parameters, inputs and outputs.

343

Any analytical procedure can be adapted on the platform and Galaxy can be used through the whole data life cycle ([https://rdmkit.elixir-europe.org/galaxy\\_assembly](https://rdmkit.elixir-europe.org/galaxy_assembly)). One can use off-the-shelf tools, workflows, and tutorials to design an analytical procedure, or suggest, develop, and share new workflows and tutorials, two aspects that do not require coding skills.

347

As each Galaxy tools are atomised and generalised elementary steps that can be articulated in a workflow, the Galaxy platform benefits from the same advantages as atomisation and generalisation and can help enhancing best practice application (tab. 2).

348

349

**Table 2 - Illustration of how the atomisation-generalisation framework and Galaxy implements and conforms to best practice.**

Reproducibility and transparency	Environment, software and package versions	Atomised-generalised code Can be indicated but possibly hard to manage Can also be set as an output of the analysis ( <i>e.g.</i> session info) Packages written in each coded elementary step or using a versioning system such as Conda	Galaxy Entirely packaged with Conda package manager and BioContainers Possibility to store analytical procedures as containers for persistent execution
	Inputs and parameters	One must keep track of different parametrisation and input settings at each computation Organisation of the analytical procedure reviewable by non-code developers Code developers might be able to detect errors as it is easier in shorter scripts Transparency over the development process achievable through Git	Automatically tracked and shareable with the “Galaxy history”
FAIR principles	Peer-review	Transparency over the development process achievable through Git	Reviewable “Galaxy history” and re-executable workflow Continuous peer-reviewed of tools with open-source code Transparency over the development process through Git The workflows can be reviewed by the Intergalactic Workflow Commission (IWC) for best practices
	Output provenance Findable	Can be tracked and reproduced in some cases If properly shared	Tracked with the “Galaxy history” and reproducible with workflow Web-based solution Unified system for data and software citation and attribution Tools can be made available on several servers Tools can be linked to tools registries and annotated with different ontologies Annotated workflows findable on WorkflowHub ( <a href="https://workflowhub.eu">https://workflowhub.eu</a> ) and Dockstore ( <a href="https://dockstore.org">https://dockstore.org</a> ) Free distribution of tools via the Galaxy ToolShed and workflows via WorkflowHub and Dockstore under an open-source licence
Technical and knowledge gaps	Accessible	If properly shared	Use different software, computational language and library versions on a single platform with the Conda package management system Workflows exportable in JSON and shareable through several standards ( <i>e.g.</i> Common Workflow Language; <i>Crusoe et al., 2022</i> and Research Object Crate; <i>Soiland-Reyes et al., 2022</i> )
	Interoperable	When properly generalised, different elementary steps should be useable in interaction with each other	Tools, histories and workflows are re-executable, reusable and adaptable with different analytical procedure, parametrisation and/or inputs. Open-source code can be used outside of a Galaxy server Tools interface, workflow annotations, help sections and tutorials are a valuable help
Collaboration and attribution	Reusable	Generalised elementary steps are reusable and adaptable with different analytical procedure, parametrisation and/or inputs The analytical procedure is clearer when properly atomised	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help
	Understanding	Learning the analytical procedure design separately from computing languages, giving structure to trainees Reusability of elementary steps for trainees Need for a computation cluster if large data or demanding algorithm Achievable through collaborative code-editing applications	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help
Collaboration and attribution	Teaching opportunities	Learning the analytical procedure design separately from computing languages, giving structure to trainees Reusability of elementary steps for trainees	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help
	Computing capacity	Need for a computation cluster if large data or demanding algorithm Achievable through collaborative code-editing applications	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help
Collaboration and attribution	Analysis design and development	Easy reuse of openly shared elementary steps could lead to higher citation rates	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help
	Citation	Easy reuse of openly shared elementary steps could lead to higher citation rates	Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help Tools interface, workflow annotations, help sections and tutorials are a valuable help

352 The Galaxy platform emphasises (i) accessibility of tools and data even without  
353 programming experience, (ii) reproducibility through the easy creation and reuse of analysis  
354 workflows, (iii) transparency through the open-source distribution of underlying codes; and  
355 (iv) community support.

356 For scientists, from a user's point of view, it offers extensive computing power and a  
357 graphical interface to use analysis workflows, even without experience in software  
358 development. Web-based access allows easy sharing of analytical workflows between  
359 collaborators and with a broader audience. Galaxy supports tools in almost any  
360 computational language, including R and Python, two of the most used languages in ecology,  
361 with many packages dedicated to ecological and biodiversity-oriented analyses incorporated  
362 (Lai *et al.*, 2019).

363 Anyone can use the tools on Galaxy and/or develop new tools and workflows to make  
364 them available to all by publishing them in the shared Galaxy ToolShed  
365 (<https://toolshed.g2.bx.psu.edu/>) which ensures that the tools and dependencies can be  
366 installed on any Galaxy servers. Any analytical procedure or workflow can be shared and  
367 enriched in parallel by several users, facilitating teamwork.

368 The platform is community-driven which permits continuous peer review of the platform  
369 and of the tools, workflows and tutorials provided. Many tutorials are available on the Galaxy  
370 Training Network (GTN; <https://training.galaxyproject.org/>) which is a valuable asset to the  
371 accessibility and reusability of tools and workflows (Batut *et al.*, 2018; Hiltmann *et al.*, 2023).

372 If enough researchers and experts start using and contributing to the platform, the number  
373 and content of available analytical procedures could expand at the same pace as latest  
374 analytical methodologies are integrated to research processes. If a different platform fits best  
375 and is more widely used by ecological and biodiversity scientific communities in the end, the  
376 work done on Galaxy will not be lost as tools are easily transposable to other interfaces (*e.g.*  
377 scripts directly usable with R, Python, etc., translation of workflows to other workflow  
378 engines).

379 Galaxy is ready to use and has proved its efficiency and suitability in other research fields,  
380 including genomics and climate science (Knijn *et al.* 2020; Serrano-Solano *et al.*, 2022).  
381 Galaxy-Ecology has implemented workflows for biodiversity data exploration, eDNA  
382 processing, general population and community metrics and models, ecoregionalisation, NDVI  
383 (Normalised difference vegetation index) computation with Sentinel-2 data among others  
384 (see some examples: <https://workflowhub.eu/workflows/657>) and tutorials for several of them  
385 are available on the GTN platform (see [https://training.galaxyproject.org/training-  
386 material/topics/ecology](https://training.galaxyproject.org/training-material/topics/ecology)).

387 In addition to using existing tools, users may develop and upload entirely new tools and  
388 workflows to the Galaxy server in any computational language to make them accessible to all  
389 other users.

390 Galaxy is a participative platform and several ways to participate to Galaxy exist depending  
391 on one's skills, available time, and needs. Anyone can participate to the Galaxy-Ecology  
392 initiative by:

- 393 – Sharing datasets, histories and workflows;
- 394 – Giving feedback on servers, tools, and workflows;
- 395 – Sharing tools and workflows ideas (eventually with code) through Git issues;

- 396 – Asking for tool modifications through issues;
- 397 – Modifying existing tools or proposing new tools through GitHub or GitLab;
- 398 – Writing or contributing to a GTN tutorial on a specific functionality or a workflow on
- 399 the Galaxy Training Network platform;
- 400 – Create learning pathways, a set of tutorials curated by community experts to form a
- 401 coherent set of lessons around a topic, building up knowledge
- 402 (<https://training.galaxyproject.org/training-material/learning-pathways>);
- 403 – Propose training events and help users in the utilisation of a workflow and tutorial.

404

405 Analyses are rarely computed only once. Any analysis with a generalisation potential is a  
406 suitable candidate to be Galaxy-fied. A methodological framework is presented in online  
407 supplementary material  
408 ([https://github.com/ColineRoyaux/Galaxy\\_Templates/blob/main/Methods/Methods%20-  
409 %20How%20to%20Galaxy-fy%20your%20analytical%20procedure\\_.md](https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure_.md)) at three levels  
410 depending on potential interests, computing language skills, and willingness to invest more or  
411 less time in the process: (i) ‘user’ relying on existing Galaxy tools and workflows to analyse  
412 data (lower time investment), (ii) ‘developer’ relying on existing and validated analytical  
413 procedure to develop Galaxy tools and workflows (highest time investment), and (iii) ‘trainer’  
414 relying on existing Galaxy tools to share workflows and create training material (variable time  
415 investment).

416

## Discussion and limitations

417 There are many best practices and recommendations existing for analytical procedures,  
418 data management, and computational code development. The levels of application of these  
419 best practices fall within a continuum offering a range of possibilities from the sole sharing of  
420 processed and interpreted results with a brief description of methods to an executable paper  
421 published within a container and emulated virtual machine (Strijkers *et al.*, 2011; Grüning *et al.*,  
422 2018). Situated somewhere in between the aforementioned extremes, the atomisation –  
423 generalisation framework and the utilisation of the Galaxy platform might represent viable  
424 solutions offering a satisfactory level of best practices.

425

426 Atomisation and generalisation of computer codes can represent a relatively low  
427 investment strategy to attain certain levels of best practices such as transparency and  
428 reusability. It also carries advantages such as easier peer review, modularity of analytical  
429 procedures and, consequently, time savings. Indeed, applying the framework is not sufficient  
430 to attain the highest levels of best practices. For reproducibility and transparency, the  
431 management of the environment, software and package versions can be hard to maintain and  
432 record. For example, on a local computer a comprehensive tracking of input, outputs and  
433 codes requires meticulous management of folder structure in the environment. Additionally,  
434 non-code developers will be able to partially review the analytical procedure only if the  
435 workflow is clearly outlined in an adapted format (*e.g.* table, graphical representation).  
436 Accessibility and findability of the atomised and generalised analytical procedure is  
437 dependent of its proper sharing (*e.g.* persistent link, open repository).



438 Galaxy can represent an easier gateway towards higher levels of best practice as sharing a  
439 complete, detailed and (re-)executable analytical procedure is facilitated through provenance  
440 tracking and automatic metadata enrichment. In comparison, many scientific workflow  
441 management systems, such as Snakemake, Nextflow or the R package Targets, operate from  
442 the command line. In ecology, numerous initiatives have tried to introduce such systems,  
443 starting with more user-friendly solutions. For example, the KNIME and Kepler systems with  
444 the CoESRA initiative (Collaborative Environment for Scholarly Research and Analysis) in  
445 Australia; Taverna with the BioVeL initiative (Biodiversity Virtual e-Laboratory) in Europe; or  
446 very recently, the BON in a Box pipeline engine. These systems are more accessible to new  
447 users by offering a graphical interface while achieving high specificity (Berthold *et al.*, 2007;  
448 Hardisty *et al.*, 2016; <https://boninabox.geobon.org/>). However, good computer programming  
449 or scientific workflow management knowledge is still necessary to use these applications  
450 appropriately.

451 In comparison to the atomisation-generalisation framework, Galaxy can be rightfully seen  
452 as necessitating more time investment for scientists with programming experience as it  
453 requires to learn to use a new platform. Additionally, more effort may be required on Galaxy  
454 when an additional analytical step needs to be developed, but the Galaxy community can be  
455 an efficient crutch on which hard-pressed scientists can rely. Indeed, one can ask for help on  
456 the implementation of tools whether one knows computing languages and can share their  
457 code or not.

458 This note showcases a simple proposition to achieve best practices in analytical  
459 procedures with two plain guidelines: atomisation and generalisation. This straightforward  
460 framework represents a different manner to think and build analytical procedures; it doesn't  
461 require using a new technology or learning to use a new software. In terms of attaining higher  
462 levels of best practice, whether it is through the atomisation-generalisation framework,  
463 Galaxy, a combination of the two or otherwise, the optimal approach is to be determined by  
464 individuals depending on their interests, projects, and available resources. Relying on existing  
465 solutions as much as possible is, in our perspective, an efficient way to achieve a better  
466 understanding of best practices and their implications. Given the current environmental crisis,  
467 science has the major political and social responsibility to maintain good levels of  
468 transparency, reproducibility and efficiency.

469

## Acknowledgements

470 Authors want to thank Sandrine Pavoine for its highly relevant and helpful advice and  
471 reviews on both the content and the form of the article. Authors are thankful to Thimothée  
472 Poisot (recommender), Nick Isaac (reviewer) and one anonymous reviewer for their advice  
473 during the Peer Community In review. Their help and suggestions on the structure and the  
474 content of the manuscript really helped to get the message of the article across in a more  
475 accessible manner.

## 476 Authors contribution statement

477 C. R. drafted the article text, tables, and figures.

478 C. R. conceptualised the atomisation – generalisation framework with J.-B. M. and Y. L.B.  
479 while working on the development of Galaxy workflows.

480 J.-B. M. and Y. L.B. reviewed and helped rewrite many parts of the draft.  
481 Y. R. and D. P. helped inspire and were invested in the early design of the article.  
482 M. J. and P. S. tested and approved the appliance of the framework.  
483 O. N., M. J., Y. R., M. E., B. B., A. F., H. R. and S. H. highly enhanced the quality of the  
484 redaction in both form and content at several stages of the draft.  
485 H. R, S. H., B. B., A. F., and B. G. are involved in the Galaxy-E initiative and provided many  
486 advice on the redaction of the article and/or on the development of the initiative.  
487 M. E. and G. M. are involved in Antarctic-oriented Galaxy tool and workflow development  
488 coordination.  
489 C. B., R. L., A. M., Y. B., A. A., T. V. and V. C. developed scripts, tools and/or Galaxy workflows  
490 to contribute to the Galaxy-E initiative.  
491 E. A. developed R scripts and apps used to integrate R Shiny apps as Galaxy interactive  
492 tools and initiate "Research Data management Galaxy tools".  
493 E. M. and C. U. developed the first training materials for Galaxy-E.  
494 E. T. worked on the use of the first Galaxy-E analysis.  
495 M. D., G. L. and R. J. were coordinating the prefiguration of Galaxy-E through the 65 Millions  
496 d'Observateurs project.  
497 Additionally, all authors reviewed and approved the article draft.

498

## Funding

499 Funding were provided by the European Union through the Erasmus+ Gallantries project;  
500 the Agence Nationale de la Recherche through the 65 Million d'Observateurs and the IA-Biodiv  
501 projects; the French National Fund for Open Science through the OpenMetaPaper project; the  
502 European commission through the H2020 EOSC-Pillar, GAPARS projects, and Horizon Europe  
503 FAIRE EASE project; the GO FAIR initiative through the BiodiFAIRse Implementation Network;  
504 the Blue Nature Alliance; and the Antarctic and Southern Ocean Coalition. Finally, funding by  
505 the French Ministry of Higher Education and Research were provided for the "Pôle national de  
506 données de biodiversité" e-infrastructure.

507

## Conflict of interest disclosure

508 The authors declare that they comply with the PCI rule of having no financial conflicts of  
509 interest in relation to the content of the article.

510

## References

511 Araújo MB, Anderson RP, Barbosa AM, Beale CM, Dormann CF, Early R, Garcia RA, Guisan A,  
512 Maiorano L, Naimi B, O'Hara RB, Zimmermann NE, Rahbek C (2019) Standards for  
513 distribution models in biodiversity assessments. *Science Advances*, **5**, 1–12.  
514 <https://doi.org/10.1126/sciadv.aat4858>  
515 Archmiller AA, Johnson AD, Nolan J, Edwards M, Elliott LH, Ferguson JM, Iannarilli F, Vélez J,  
516 Vitense K, Johnson DH, Fieberg J (2020) Computational Reproducibility in The Wildlife  
517 Society's Flagship Journals. *Journal of Wildlife Management*, **84**, 1012–1017.  
518 <https://doi.org/10.1002/JWMG.21855>

519 Batut B, Hiltemann S, Bagnacani A, Baker D, Bhardwaj V, Blank C, Bretaudeau A, Brillet-  
520 Guéguen L, Čech M, Chilton J, Clements D, Doppelt-Azeroual O, Erxleben A, Freeberg MA,  
521 Gladman S, Hoogstrate Y, Hotz HR, Houwaart T, Jagtap P, Larivière D, Le Corguillé G, Manke  
522 T, Mareuil F, Ramírez F, Ryan D, Sigloch FC, Soranzo N, Wolff J, Videm P, Wolfien M, Wubuli  
523 A, Yusuf D, Taylor J, Backofen R, Nekrutenko A, Grüning B (2018) Community-Driven Data  
524 Analysis Training for Biology. *Cell Systems*, **6**, 752-758.  
525 <https://doi.org/10.1016/j.cels.2018.05.012>

526 Berthold MR, Cebron N, Dill F, Gabriel TR, Kötter T, Meinel T, Ohl P, Sieb C, Thiel K, Wiswedel B  
527 (2007) KNIME: The Konstanz Information Miner. *Studies in Classification, Data Analysis, and*  
528 *Knowledge Organization*, 319–326. [https://doi.org/10.1007/978-3-540-78246-9\\_38](https://doi.org/10.1007/978-3-540-78246-9_38)

529 Borgman CL (2020) Qu'est-ce que le travail scientifique des données? Big data, little data, no  
530 data. <https://doi.org/10.4000/BOOKS.OEP.14692>

531 Boyd RJ, August TA, Cooke R, Logie M, Mancini F, Powney GD, Roy DB, Turvey K, Isaac NJB  
532 (2023) An operational workflow for producing periodic estimates of species occupancy at  
533 national scales. *Biological Reviews*, 98, 1492–1508. <https://doi.org/10.1111/brv.12961>

534 Carroll S, Garba I, Figueroa-Rodríguez O, Holbrook J, Lovett R, Materechera S, Parsons M,  
535 Raseroka K, Rodriguez-Lonebear D, Rowe R, Sara R, Walker J, Anderson J, Hudson M (2020)  
536 The CARE Principles for Indigenous Data Governance. *Data Science Journal*, **19**, 43.  
537 <https://doi.org/10.5334/dsj-2020-043>

538 Casajus N. (2023) rcompendium: An R package to create a package or research compendium  
539 structure. <https://github.com/FRBCesab/rcompendium>

540 Cohen-Boulakia S, Belhajjame K, Collin O, Chopard J, Froidevaux C, Gaignard A, Hinsén K,  
541 Larmande P, Bras Y Le, Lemoine F, Mareuil F, Ménager H, Pradal C, Blanchet C (2017)  
542 Scientific workflows for computational reproducibility in the life sciences: Status,  
543 challenges and opportunities. *Future Generation Computer Systems*, **75**, 284–298.  
544 <https://doi.org/10.1016/j.future.2017.01.012>

545 Crusoe MR, Abeln S, Iosup A, Amstutz P, Chilton J, Tijanić N, Ménager H, Soiland-Reyes S,  
546 Goble C (2022) Methods Included: Standardizing Computational Reuse and Portability with  
547 the Common Workflow Language. *Communications of the ACM*, **65**, 54–63.  
548 <https://doi.org/10.1145/3486897>

549 Culina A, van den Berg I, Evans S, Sánchez-Tójar A (2020) Low availability of code in ecology: A  
550 call for urgent action. *PLOS Biology*, **18**, e3000763.  
551 <https://doi.org/10.1371/JOURNAL.PBIO.3000763>

552 Di Cosmo R, Zacchiroli S (2017) Software Heritage: Why and How to Preserve Software Source  
553 Code. <https://hal.science/hal-01590958>

554 Di Tommaso P, Chatzou M, Floden EW, Barja P., Palumbo E, Notredame C (2017) Nextflow  
555 enables reproducible computational workflows. *Nature Biotechnology*, **35**, 316–319.  
556 <https://doi.org/10.1038/nbt.3820>

557 Ellemers N (2021) Science as collaborative knowledge generation. *British Journal of Social*  
558 *Psychology*, **60**, 1–28. <https://doi.org/10.1111/BJSO.12430>

559 EMBL Australia Bioinformatics Resource (2013) Community Survey Report [https://www.embl-](https://www.embl-abr.org.au/news/braembl-community-survey-report-2013/)  
560 [abr.org.au/news/braembl-community-survey-report-2013/](https://www.embl-abr.org.au/news/braembl-community-survey-report-2013/)

561 Emery NC, Crispo E, Supp SR, Farrell KJ, Kerkhoff AJ, Bledsoe EK, O'Donnell KL, McCall AC,  
562 Aiello-Lammens ME (2021) Data Science in Undergraduate Life Science Education: A Need

563 for Instructor Skills Training. *BioScience*, **71**, 1274–1287.  
564 <https://doi.org/10.1093/BIOSCI/BIAB107>

565 European Commission, Directorate-General for Research and Innovation (2018) Cost-benefit  
566 analysis for FAIR research data: cost of not having FAIR research data. *Publications Office*.  
567 <https://doi.org/10.2777/02999>

568 Fanelli D (2018) Is science really facing a reproducibility crisis, and do we need it to?  
569 *Proceedings of the National Academy of Sciences of the United States of America*, **115**,  
570 2628–2631. <https://doi.org/10.1073/pnas.1708272114>

571 Fang FC, Casadevall A (2015) Competitive Science: Is Competition Ruining Science? *Infection*  
572 *and Immunity*, **83**, 1229–1233. <https://doi.org/10.1128/IAI.02939-14>

573 Farley SS, Dawson A, Goring SJ, Williams JW (2018) Situating Ecology as a Big-Data Science:  
574 Current Advances, Challenges, and Solutions. *BioScience*, **68**, 563–576.  
575 <https://doi.org/10.1093/BIOSCI/BIY068>

576 Field B, Booth A, Illott I, Gerrish K (2014) *Using the Knowledge to Action Framework in practice:*  
577 *a citation analysis and systematic review. Implementation Science*, **9**, 172.  
578 <https://doi.org/10.1186/s13012-014-0172-2>

579 Goble C, Cohen-Boulakia S, Soiland-Reyes S, Garijo D, Gil Y, Crusoe MR, Peters K, Schober D  
580 (2020) FAIR Computational Workflows. *Data Intelligence*, **2**, 108–121.  
581 [https://doi.org/10.1162/dint\\_a\\_00033](https://doi.org/10.1162/dint_a_00033)

582 Gomes DGE, Pottier P, Crystal-Ornelas R, Hudgins EJ, Foroughirad V, Sánchez-Reyes LL, Turba  
583 R, Martinez PA, Moreau D, Bertram MG, Smout CA, Gaynor KM (2022) Why don't we share  
584 data and code? Perceived barriers and benefits to public archiving practices. *Proceedings*  
585 *of the Royal Society B*, **289**, 20221113 <https://doi.org/10.1098/rspb.2022.1113>

586 Gownaris NJ, Vermeir K, Bittner MI, Gunawardena L, Kaur-Ghumaan S, Lepenies R, Ntsefong  
587 GN, Zakari IS (2022) Barriers to Full Participation in the Open Science Life Cycle among  
588 Early Career Researchers. *Data Science Journal*, **21**, 2. [https://doi.org/10.5334/DSJ-2022-](https://doi.org/10.5334/DSJ-2022-002)  
589 [002](https://doi.org/10.5334/DSJ-2022-002)

590 Grüning B, Chilton J, Köster J, Dale R, Soranzo N, van den Beek M, Goecks J, Backofen R,  
591 Nekrutenko A, Taylor J (2018) Practical Computational Reproducibility in the Life Sciences.  
592 *Cell Systems*, **6**, 631–635. <https://doi.org/10.1016/j.cels.2018.03.014>

593 Hampton SE, Jones MB, Wasser LA, Schildhauer MP, Supp SR, Brun J, Hernandez RR, Boettiger  
594 C, Collins SL, Gross LJ, Fernández DS, Budden A, White EP, Teal TK, Labou SG, Aukema JE  
595 (2017) Skills and Knowledge for Data-Intensive Environmental Research. *BioScience*, **67**,  
596 546–557. <https://doi.org/10.1093/BIOSCI/BIX025>

597 Hardisty AR, Bacall F, Beard N, Balcázar-Vargas MP, Balech B, Barcza Z, Bourlat SJ, Giovanni R,  
598 Jong Y, Leo F, Dobor L, Donvito G, Fellows D, Guerra AF, Ferreira N, Fetyukova Y, Fosso B,  
599 Giddy J, Goble C, Güntsch A, Haines R, Ernst VH, Hettling H, Hidy D, Horváth F, Ittész D,  
600 Ittész P, Jones A, Kottmann R, Kulawik R, Leidenberger S, Lyytikäinen-Saarenmaa P,  
601 Mathew C, Morrison N, Nenadic A, Hidalgo AN, Obst M, Oostermeijer G, Paymal E, Pesole G,  
602 Pinto S, Poigné A, Fernandez FQ, Santamaria M, Saarenmaa H, Sipos G, Sylla KH, Tähtinen  
603 M, Vicario S, Vos RA, Williams AR, Yilmaz P (2016) BioVeL: A virtual laboratory for data  
604 analysis and modelling in biodiversity science and ecology. *BMC Ecology*, **16**, 49.  
605 <https://doi.org/10.1186/S12898-016-0103-Y>

606 Hiltemann S, Rasche H, Gladman S, Hotz HR, Larivière D, Blankenberg D, Jagtap PD, Wollmann  
607 T, Bretaudeau A, Goué N, Griffin TJ, Royaux C, Bras Y Le, Mehta S, Syme A, Coppens F,  
608 Droesbeke B, Soranzo N, Bacon W, Psomopoulos F, Gallardo-Alba C, Davis J, Föll MC,  
609 Fahrner M, Doyle MA, Serrano-Solano B, Fouilloux AC, van Heusden P, Maier W, Clements D,  
610 Heyl F, Grüning B, Batut B (2023) Galaxy Training: A powerful framework for teaching! *PLOS*  
611 *Computational Biology*, **19**, e1010752. <https://doi.org/10.1371/JOURNAL.PCBI.1010752>

612 Ioannidis JPA (2022) Correction: Why Most Published Research Findings Are False. *Plos*  
613 *Medicine*, **39**, e1004085. <https://doi.org/10.1371/JOURNAL.PMED.1004085>

614 Ivimey-Cook ER, Pick JL, Bairos-Novak K, Culina A, Gould E, Grainger M, Marshall B, Moreau D,  
615 Paquet M, Royauté R, Sanchez-Tojar A, Silva I, Windecker S (2023) Implementing Code  
616 Review in the Scientific Workflow: Insights from Ecology and Evolutionary Biology.  
617 *EcoEvoRxiv*. <https://doi.org/10.32942/X2CG64>

618 Jenkins GB, Beckerman AP, Bellard C, Benítez-López A, Ellison AM, Foote CG, Hufton AL,  
619 Lashley MA, Lortie CJ, Ma Z, Moore AJ, Narum SR, Nilsson J, O’Boyle B, Provete DB, Razgour  
620 O, Rieseberg L, Riginos C, Santini L, Sibbett B, Peres-Neto PR (2023) Reproducibility in  
621 ecology and evolution: Minimum standards for data and code. *Ecology and Evolution*, **13**,  
622 e9961. <https://doi.org/10.1002/ECE3.9961>

623 Jetz W, McGeoch MA, Guralnick R, Ferrier S, Beck J, Costello MJ, Fernandez M, Geller GN, Keil P,  
624 Merow C, Meyer C, Muller-Karger FE, Pereira HM, Regan EC, Schmeller DS, Turak E (2019)  
625 Essential biodiversity variables for mapping and monitoring species populations. *Nature*  
626 *Ecology and Evolution*, **3**, 539–551. <https://doi.org/10.1038/s41559-019-0826-1>

627 Keenan M, Cutler P, Marks J, Meylan R, Smith C, Koivisto E (2012) Orienting international  
628 science cooperation to meet global “grand challenges.” *Science and Public Policy*, **39**, 166–  
629 177. <https://doi.org/10.1093/SCIPOL/SCS019>

630 Knijn A, Michelacci V, Orsini M, Morabito S (2020) Advanced Research Infrastructure for  
631 Experimentation in genomicS (ARIES): a lustrum of Galaxy experience. *bioRxiv*.  
632 <https://doi.org/10.1101/2020.05.14.095901>

633 Köster J, Rahmann S (2012) Snakemake—a scalable bioinformatics workflow engine.  
634 *Bioinformatics*, **28**, 2520–2522. <https://doi.org/10.1093/bioinformatics/bts480>

635 Lai J, Lortie CJ, Muenchen RA, Yang J, Ma K (2019) Evaluating the popularity of R in ecology.  
636 *Ecosphere*, **10**, e02567. <https://doi.org/10.1002/ECS2.2567>

637 Lamprecht A-L, Garcia L, Kuzak M, Martinez C, Arcila R, Martin Del Pico E, Dominguez Del Angel  
638 V, van de Sandt S, Ison J, Martinez PA, McQuilton P, Valencia A, Harrow J, Psomopoulos F,  
639 Gelpi JL, Chue Hong N, Goble C, Capella-Gutierrez S (2019) Towards FAIR principles  
640 for research software. *Data Science*, **3**, 37–59. <https://doi.org/10.3233/ds-190026>

641 Larcombe L, Hendricusdottir R, Attwood T, Bacall F, Beard N, Bellis L, Dunn W, Hancock J,  
642 Nenadic A, Orengo C, Overduin B, Sansone S, Thurston M, Viant M, Winder C, Goble C,  
643 Ponting C, Rustici G (2017) ELIXIR-UK role in bioinformatics training at the national level  
644 and across ELIXIR. *F1000Research*, **6**, 952. <https://doi.org/10.12688/f1000research.11837.1>

645 Leroy B (2023) Choosing presence-only species distribution models. *Journal of Biogeography*,  
646 **50**, 247–250. <https://doi.org/10.1111/jbi.14505>

647 Lortie CJ (2021) The early bird gets the return: The benefits of publishing your data sooner.  
648 *Ecology and Evolution*, **11**, 10736–10740. <https://doi.org/10.1002/ECE3.7853>



649 McIntire EJB, Chubaty AM, Cumming SG, Andison D, Barros C, Boisvenue C, Haché S, Luo Y,  
650 Micheletti T, Stewart FEC (2022) PERFICT: A Re-imagined foundation for predictive ecology.  
651 *Ecology Letters*, **25**, 1345–1351. <https://doi.org/10.1111/ELE.13994>

652 Michener WK (2015) Ten Simple Rules for Creating a Good Data Management Plan. *PLOS*  
653 *Computational Biology*, **11**, e1004525. <https://doi.org/10.1371/JOURNAL.PCBI.1004525>

654 Michener WK, Jones MB (2012) Ecoinformatics: Supporting ecology as a data-intensive science.  
655 *Trends in Ecology and Evolution*, **27**, 85–93. <https://doi.org/10.1016/j.tree.2011.11.016>

656 Minocher R, Atmaca S, Bavero C, McElreath R, Beheim B (2021) Estimating the reproducibility  
657 of social learning research published between 1955 and 2018. *Royal Society Open Science*,  
658 **8**, 210450. <https://doi.org/10.1098/RSOS.210450>

659 Munafò MR, Nosek BA, Bishop DVM, Button KS, Chambers CD, Percie Du Sert N, Simonsohn U,  
660 Wagenmakers EJ, Ware JJ, Ioannidis JPA (2017) A manifesto for reproducible science.  
661 *Nature Human Behaviour*, **1**, 0021. <https://doi.org/10.1038/s41562-016-0021>

662 Natural Environment Research Council (2010, 2012) Most Wanted: Postgraduate Skills Needs  
663 in the Environment Sector.

664 Plesser HE (2018) Reproducibility vs. Replicability: A brief history of a confused terminology.  
665 *Frontiers in Neuroinformatics*, **11**, 76. <https://doi.org/10.3389/FNINF.2017.00076>

666 Powers SM, Hampton SE (2019) Open science, reproducibility, and transparency in ecology.  
667 *Ecological applications*, **29**, e01822. <https://doi.org/10.1002/eap.1822>

668 Serrano-Solano B, Fouilloux A, Eguinoa I, Kalaš M, Grüning B, Coppens F (2022) Galaxy: A  
669 Decade of Realising CWFR Concepts. *Data Intelligence*, **4**, 358–371.  
670 [https://doi.org/10.1162/dint\\_a\\_00136](https://doi.org/10.1162/dint_a_00136)

671 Soiland-Reyes S, Sefton P, Crosas M, Castro LJ, Coppens F, Fernández JM, Garijo D, Grüning B,  
672 La Rosa M, Leo S, Ó Carragáin E, Portier M, Trisovic A, Community R-C, Groth P, Goble C  
673 (2022) Packaging research artefacts with RO-Crate. *Data Science*, **5**, 97–138.  
674 <https://doi.org/10.3233/DS-210053>

675 Strijkers R, Cushing R, Vasyunin D, De Laat C, Belloum ASZ, Meijer R (2011) Toward executable  
676 scientific publications. *Procedia Computer Science*, **4**, 707–715.  
677 <https://doi.org/10.1016/J.PROCS.2011.04.074>

678 The Galaxy Community (2022) The Galaxy platform for accessible, reproducible and  
679 collaborative biomedical analyses: 2022 update. *Nucleic acids research*, **50**, W345–W351.  
680 <https://doi.org/10.1093/NAR/GKAC247>

681 Touchon JC, McCoy MW (2016) The mismatch between current statistical practice and  
682 doctoral training in ecology. *Ecosphere*, **7**, e01394. <https://doi.org/10.1002/ECS2.1394>

683 Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten  
684 JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M, Dillo I,  
685 Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, Gray AJG, Groth P, Goble C,  
686 Grethe JS, Heringa J, t Hoen PAC, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ, Martone ME,  
687 Mons A, Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R, Sansone SA, Schultes E,  
688 Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, Van Der Lei J, Van Mulligen E,  
689 Velterop J, Waagmeester A, Wittenburg P, Wolstencroft K, Zhao J, Mons B (2016) Comment:  
690 The FAIR Guiding Principles for scientific data management and stewardship. *Scientific*  
691 *Data*, **3**, 1–9. <https://doi.org/10.1038/sdata.2016.18>

692 Williams JJ, Teal TK (2017) A vision for collaborative training infrastructure for bioinformatics.  
693 *Annals of the New York Academy of Sciences*, **1387**, 54–60.  
694 <https://doi.org/10.1111/NYAS.13207>

695 Zurell D, Franklin J, König C, Bouchet PJ, Dormann CF, Elith J, Fandos G, Feng X, Guillera-  
696 Arroita G, Guisan A, Lahoz-Monfort JJ, Leitão PJ, Park DS, Peterson AT, Rapacciuolo G,  
697 Schmatz DR, Schröder B, Serra-Diaz JM, Thuiller W, Yates KL, Zimmermann NE, Merow C  
698 (2020) A standard protocol for reporting species distribution models. *Ecography*, **43**, 1261–  
699 1277. <https://doi.org/10.1111/ecog.04960>