

# 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>, Segueineau 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>, Lois 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

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

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

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

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

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

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

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

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

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

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

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

66

67 \*Corresponding author

68 Correspondence: coline.royaux@mnhn.fr

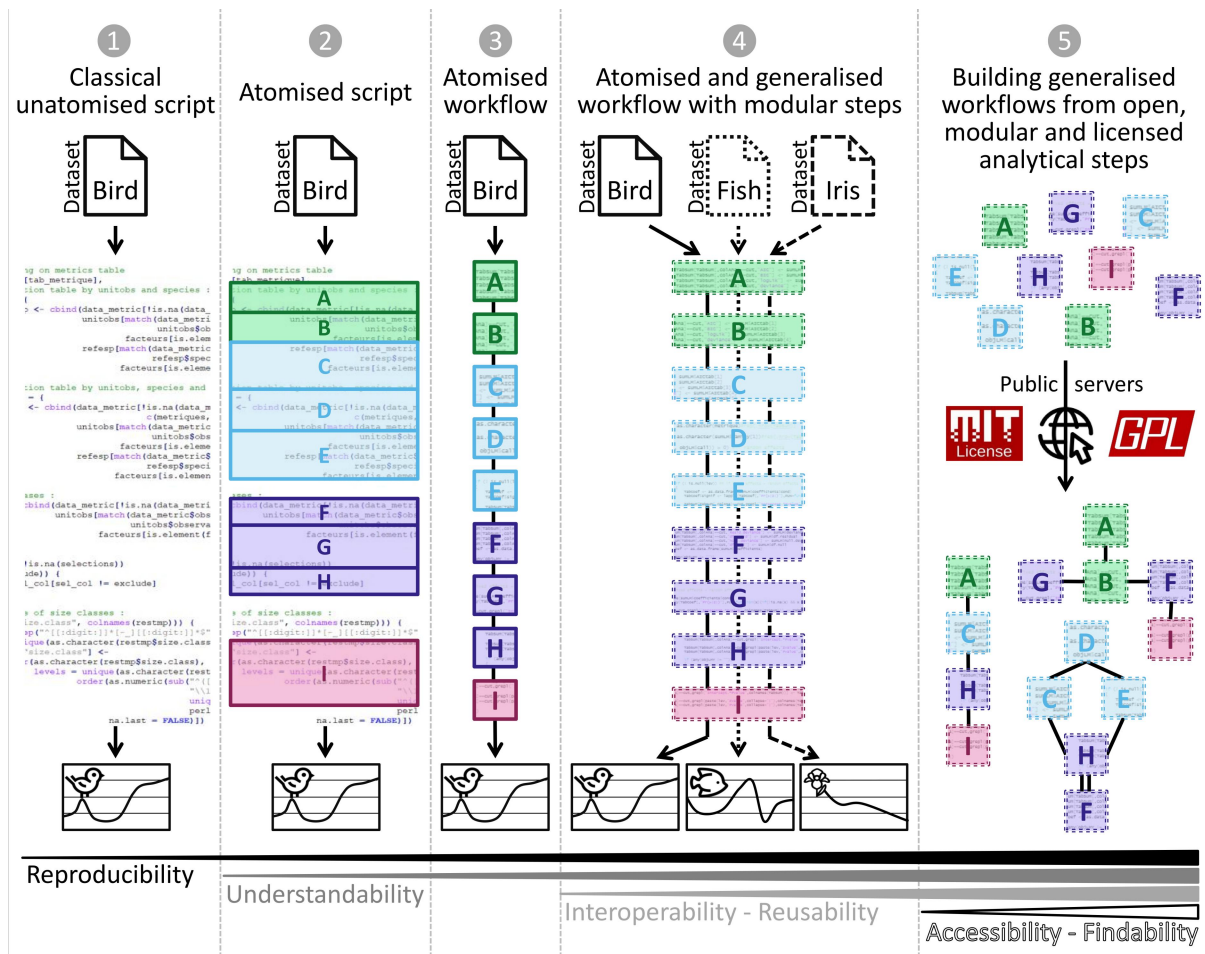
69

## 70 **ABSTRACT**

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

88

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



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

## 96 Ecology's Reproducibility Crisis

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

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

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

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


138 Although data accessibility has been substantially improved in ecology  
139 during the past decade, sharing analytical scripts and codes remain largely

140 marginal (Archmiller *et al.*, 2020; Culina *et al.*, 2020; Minocher *et al.*, 2021;  
141 Ivimey-Cook *et al.*, 2023).

142 However, even if sharing code is necessary to achieve good computational  
143 reproducibility, it is insufficient. Therefore, the utilisation of computational  
144 workflows has been suggested as a solution for improving computational  
145 reproducibility (Cohen-Boulakia *et al.*, 2017; Grüning *et al.*, 2018) through  
146 software such as Snakemake (Köster & Rahmann, 2012), Nextflow (Di  
147 Tommaso *et al.*, 2017), or Galaxy (The Galaxy Community, 2022). A workflow  
148 is generally defined as a sequence of distinct computational tasks for a  
149 particular objective (Goble *et al.*, 2020). As such, a workflow represents the  
150 backbone of a single specific analysis. Throughout the analytical procedure, a  
151 typical workflow starts with raw data, which can be extracted from several  
152 databases or data files and processed through a series of analytical steps.  
153 The products resulting from these analytical steps (*i.e.* the outputs of the  
154 computational workflow) can be data files, graphic representations and any  
155 associated metrics.

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

- 159 – the specificities of the workflow, the analysis steps and associated  
160 tools;
- 161 – the workflow entries, datasets and parameters;
- 162 – the environment and context of the use of the workflow;
- 163 – the results obtained and the outputs of the workflow.

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

## 180 Impact on Ecology Research

181 The efficiency of the expertise and research is greatly affected by the lack  
182 of computational reproducibility and FAIRness of analytical procedures. FAIR  
183 research data was estimated to save 10.2 billion € per year in Europe  
184 (Munafò *et al.*, 2017; European commission, 2018; Gomes *et al.*, 2022).  
185 Moreover, consistent application of reproducibility and FAIR principles will

186 improve trust in research studies and scientific reports (Powers & Hampton,  
187 2019; Lortie, 2021; Jenkins *et al.*, 2023).

188 The widespread use of computational languages to process large-scale  
189 data and analyse complex systems has been a major advance in studying the  
190 ecosphere at any spatio-temporal scale (Michener & Jones, 2012; Farley *et al.*,  
191 2018). However, the ever-growing technical and programming skills required  
192 to take advantage of such computational solutions by the scientific  
193 community raise new challenges (Jetz *et al.*, 2019; Leroy, 2022; Boyd *et al.*,  
194 2023). The use of increasingly complex analytical solutions, paired with  
195 different approaches or programming languages, mechanically reduces the  
196 number of potential users, limiting collaboration and fragilising fundamental  
197 pillars of scientific knowledge such as the peer-review process and critical  
198 evaluation. As a response to this situation, adequate training was identified  
199 by life science researchers (*Community Survey Report*, 2013; Williams & Teal,  
200 2017; Larcombe *et al.*, 2017), as it would help involve more people in the  
201 understanding of current analytical solutions and benefit to scientific  
202 cooperation (Touchon & McCoy, 2016; Gownaris *et al.*, 2022). Research is  
203 typically structured through a highly competitive organisation, with a  
204 potentially detrimental effect on scientific knowledge (Fang & Casadevall,  
205 2015). Instead, fostering collaboration and collective intelligence by  
206 promoting transparent sharing of analytical procedures, would offer more  
207 persistent and robust ways to achieve actionable science (Ellemers, 2021).  
208 Such efforts would be of paramount importance in environmental sciences  
209 and the conservation of biodiversity by providing governance and guiding  
210 actions with increasingly robust evidence (Keenan *et al.*, 2012).

211 Are there simple and ready-to-use solutions?

212 In this note, we aim to promote the reuse of existing concepts and  
213 solutions as pillars toward better practices for ecological analyses by  
214 providing a streamlined framework. We believe the atomisation-  
215 generalisation framework presented in the second part of this note  
216 represents an operational and actionable path for researchers and experts to  
217 attain levels of best practices (e.g. reproducibility, FAIR, open science, R  
218 compendium; Casajus N., 2023) with no more investment than they are able  
219 or willing to provide (Field *et al.*, 2014). Atomisation is used to refer to the  
220 identification of single analytical steps constituting an analytical procedure. It  
221 is a non-standard term introduced in this note to convey the idea of analytical  
222 “atoms”. As for atom particles that etymologically correspond to “indivisible”  
223 but are composed of subatomic particles, an analytical atom represents a  
224 single analytical step composed of several functions. Generalisation involves  
225 the alteration of an analytical step to enlarge its applicability in diverse  
226 contexts and for diverse purposes.

227 This framework has been formalised while building the Galaxy-Ecology  
228 (Galaxy-E) initiative (see section III). Galaxy (The Galaxy Community, 2022) is  
229 a workflow-oriented web platform for sharing and processing data. It allows  
230 scientists to share, develop, and use various datasets and data processing  
231 tools (e.g. data formatting, statistical tests, graphic representations).

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

239 Galaxy-E is a demonstration platform for applying best practices such as  
240 the FAIR principles and computational reproducibility for analytical  
241 procedures in ecology. Hence, this technical note is partly Galaxy-oriented,  
242 not to present the platform as a prescriptive solution but to give an  
243 operational example of the best practices it helps to achieve.

## 244 Framework towards best practices

### 245 Atomisation: what is it and why?

246 Atomisation refers to dividing an analytical procedure into several specific  
247 steps (“atoms”; Graphical abstract ②) generating a suite of elementary  
248 analytical steps as pictured in the Graphical abstract ③. Breaking down the  
249 analytical process into atoms functioning as building blocks allows for better  
250 understanding, modularity, and visibility of the analytical flow. It permits  
251 making it more accessible to a broader audience or facilitating the peer-  
252 review process. Indeed, an extended one-block code that imports raw data,  
253 makes pre-processing steps (e.g. filter, formatting), conducts analyses (e.g.  
254 distribution study, modelling), and performs final representations of results  
255 (e.g. maps, plots) can be challenging to understand and reuse by others or  
256 even the same person after some time.

257 McIntire *et al.* (2022) described the PERFICT approach (Prediction,  
258 Evaluation, Reusability, Free access, Interoperability, Continuous workflows,  
259 and routine Tests) to set a new foundation for models in predictive ecology.  
260 This can be applied more generally to the analytical procedure in ecology and  
261 biodiversity. In their article, McIntire and collaborators make an analogy  
262 between code development and Lego® construction, similar to our definition  
263 of atomisation. Functions are a workflow’s most fundamental analytical steps  
264 and can be seen as modular pieces, alike single pieces of Lego®. Modules  
265 can be created from a single or series of successive functions comparably as  
266 in Lego® structures made of several pieces (e.g. meant to build cars, houses,  
267 or road). These modules (or atoms, tools) can be used as standalone or  
268 combined to make simple to complex analytical workflows (e.g. data  
269 formatting or curation, running statistical models, or generating graphical  
270 elements for visualisation). Doing so, the atomisation approach may facilitate  
271 sharing or teaching analytical practices since beginners can easily  
272 understand the general organisation of the analytical procedure by simply  
273 reading the list of steps in the analysis with a limited degree of complexity.  
274 Decoupling programming skills from analytical skills can make data  
275 processing more accessible to a wider audience. Indeed, once each  
276 elementary step is clearly identified and delimited along the atomisation  
277 process, it is easier to grasp the whole analytical procedure and focus on the

278 review of each step at a time or (re)use it. New workflows can further be  
279 generated by recombining existing, validated or peer-reviewed elementary  
280 steps in innovative ways. This process can save time, increase confidence,  
281 and avoid potential programming mistakes, allowing greater focus on  
282 understanding the analytical workflow.

283 Generalisation: what is it and why?

284 Generalisation refers to the modification of an analytical procedure to  
285 make it applicable to many settings, by removing specificities related to a  
286 particular data file or data format. Generalisation aims to optimise the  
287 reusability at different times (e.g. regular result update), enlarge the  
288 application of a given analysis to different input data files while keeping the  
289 initial analytical procedure fully reproducible as pictured in the Graphical  
290 abstract ④. Generalising an analytical step requires identifying key steps and  
291 invariant parameters from those that must be adaptable to allow for the  
292 analysis to be applied to specific characteristics of various datasets. These  
293 parameters must be implemented to be easily modified if needed.  
294 Generalisation can be tricky because the higher the flexibility of an analytical  
295 step, the greater the risk of errors in its use. This is why generalisation should  
296 be complemented by clear statement and an implementation of red flags and  
297 warnings to prevent such events. As with atomisation, generalisation is  
298 primarily a conceptual way to build analytical procedures. It requires minor  
299 change of practices to reach certain degree of generalisation, avoiding  
300 additional effort later on for reusability, reproducibility, and share.

301 How to do atomisation and generalisation with computer codes: Finding  
302 balance

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



Raw script

```

# depending on metrics table
as.table(tab_metrics)
# Observation table by units and species :
nit_sp = [
  restmp <- cbind(data_metric[is.na(data_metric[, metric
  unitobs[match(data_metric$observation.unit,
  unitobs$observation.unit)]
  factors[is.element(factors, c(
  refesp$species.code)]
  refesp$species.code), colname
  factors[is.element(factors, colname
  )].
# Observation table by units, species and size class :
nit_sp_size = [
  restmp <- cbind(data_metric[is.na(data_metric[, metric
  observations[match(data_metric$observation.unit,
  observations$observation.unit)]
  factors[is.element(factors, c(
  refesp$species.code)]
  refesp$species.code), colname
  factors[is.element(factors, colname
  )].
# Other cases :
estmp <- cbind(data_metric[is.na(data_metric[, metric]
  observations[match(data_metric$observation.unit,
  observations$observation.unit)]
  factors[is.element(factors, colname
  )].
# Which (is.na(selections))
null(exclude) {
  col <- sel_col[sel_col != exclude]

  ular case of size classes :
  element("size.class", colname(restmp)) {
  length(grep("[[:digit:]]+|[[:alpha:]]+", unique(as.chu
  length(unique(as.character(restmp$size.class))) {
  restmp[, "size.class"] <-
  factor(as.character(restmp$size.class),
  levels = unique(as.character(restmp$size.class))
  order(as.numeric(sub("[[:digit:]]+", "",
  "\\.\\.",
  unique(as.character(rest
  perl = TRUE))
  na.last = FALSE))

# and density conversion -> /Um^3 :
is.element(colnames(restmp), c("biomass", "density",
  "biomass.max", "density.max",
  "biomass.sd", "density.sd"))
sp[, is.element(colnames(restmp),
  c("biomass", "density",
  "biomass.max", "density.max",
  "biomass.sd", "density.sd")) <- (sp *
  restmp[, is.element(colnames(restmp),
  c("biomass", "de
  "biomass.max",
  "biomass.sd", "den
  
```

Level 1 of atomisation

**Data exploration**

**Pre-processing**

**Analysis**  
(e. g. statistics, simulation)

Level 2 of atomisation

**Sampling plan**

**NA repartition**

**Data resolution**

**Data distribution**

**Formatting**

**Corrections**

**Filter**

**Anonymisation**

**Variable exploration**

**Tests & models**

**Projections**

**Plots**

**Maps**

Level N of atomisation  
...

318

319

Figure 1 - Illustration of the atomisation of an existing code

**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 (e.g. AIC, Jaccard) Validation methods
	Projections	Geographical projections Temporal projections
...	...	...
Representation	Plot	Raw variables Modelled results
	Map	Observations Projections
...	...	...

321 A few changes in code-writing habits can enhance the reusability of the  
322 analytical procedure by generating easy-to-understand analytical procedure  
323 without investing much time. It is best to develop each elementary step  
324 directly in separate code files and to give details of the order in which  
325 elementary steps are used for each analytical workflow. To ensure  
326 reproducibility and traceability of the results, each computation of the  
327 analytical workflow should be associated with the details of the parameters  
328 settings and datasets used. From a practical point of view, a couple of  
329 recommendations could be made for coding elementary steps in order to  
330 facilitate generalisation and ease the reuse. Once each elementary step is  
331 defined, we recommend all dependencies (e.g. software version, packages,  
332 libraries and their versions) to be set at the same place, at the start of the  
333 code, followed by modular parameters (e.g. input file location and name,  
334 column selection, modelling parameters, data specificities, output saving  
335 location). When the script of the elementary step is completed, modular  
336 parameters should be the only part of the code that may be modified in  
337 future reuse. Dependencies and subsequent computational tasks should be

338 left untouched to ensure the integrity of the analysis and then, reproducibility.  
339 In the end, it is best to add an open-source license to any analytical  
340 procedure shared publicly (e.g. MIT, GPL). It permits to clearly state the terms  
341 and conditions of diffusion, share and reuse.

342 As such, atomisation and generalisation may overcome social or  
343 psychological barriers related to transparent sharing, either related to  
344 securing ownership (e.g. DOI) and to embarrassment or fear during a peer-  
345 review process (Gomes *et al.*, 2022).

346 Atomisation and generalisation are related and complementary concepts.  
347 Atomisation into adequate elementary steps is necessary to properly  
348 generalise an analytical procedure as it permits to enhance the modularity of  
349 the procedure and its capacity to be tailored to different data types.  
350 Atomisation and generalisation must be applied from the earliest stages of  
351 the programming development of any analytical procedure in order to  
352 achieve:

- 353 – Greater transparency, even for beginners, since the relevance and  
354 coherence of each step and their successive arrangement along the  
355 analytical procedure should be appraised independently of the  
356 programming skills;
- 357 – Time savings;
- 358 – Greater reusability;
- 359 – Modularity of the elementary steps, to rearrange them differently if  
360 needed.

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

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

369 Galaxy-E is mostly aimed at scientists that process biodiversity data and  
370 already have an understanding of the general functioning of the analytical  
371 procedures they want to produce. The rationale for a user would be to create  
372 or reuse analytical workflows with high FAIRness in a collaborative and open  
373 source platform. It can be used for individual analyses as well as for  
374 collaborative projects. In some cases, if the analytical procedure is already  
375 clearly defined, it can be used by citizens or for teaching.

376 It benefits from the same advantages as the framework presented in the  
377 previous section and can help achieve a further level of FAIRness as a  
378 demonstration platform to package analyses in an accessible and user-  
379 friendly manner (tab. 2).

380  
381

**Table 2** - Comparison between the atomisation-generalisation framework and Galaxy for the achievement of best practices. Limitations are occasionally raised with short advice to mitigate them when relevant

		Atomised-generalised code	Galaxy
Reproducibility and transparency	Environment, software and package versions	Can be indicated but possibly hard to manage Can also be set as an output of the analysis (e.g. session info) Packages written in each coded elementary step or using a versioning system such as Conda	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	Automatically tracked and shareable with the “Galaxy history”
	Peer-review	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	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	Can be tracked and reproduced in some cases	Tracked with the “Galaxy history” and reproducible with workflow
FAIR principles	Findable	If properly shared	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> )
	Accessible	If properly shared	Free distribution of tools via the Galaxy ToolShed and workflows via WorkflowHub and Dockstore under an open-source licence
	Interoperable	When properly generalised, different elementary steps should be useable in interaction with each other	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 (e.g. Common Workflow Language; Crusoe <i>et al.</i> , 2022 and Research Object Crate; Soiland-Reyes <i>et al.</i> , 2022)
	Reusable	Generalised elementary steps are reusable and adaptable with different analytical procedure, parametrisation and/or inputs	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
Technical and knowledge gaps	Understandability	The analytical procedure is clearer when properly atomised	Tools interface, workflow annotations, help sections and tutorials are a valuable help
	Teaching opportunities	Learning the analytical procedure design separately from computing languages, giving structure to trainees Reusability of elementary steps for trainees	Experimenting with intricate analyses without computer code first Tutorials and videos from Galaxy Training Network ( <a href="https://training.galaxyproject.org">https://training.galaxyproject.org</a> ) Galaxy community
	Computing capacity	Need for a computation cluster if large data or demanding algorithm	HPC (High Performance Computing) through an interface Bulk (meta)data manipulation
Collaboration and attribution	Analysis design and development	Achievable through collaborative code-editing applications	With anyone through a Galaxy server
	Citation	Easy reuse of openly shared elementary steps could lead to higher citation rates	Each tool, workflow, and tutorial are provided with a unique identifier for proper attribution and citation

382

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

387 Galaxy is ready to use and has proved its efficiency and suitability in other  
388 research fields, including genomics and climate science (Knijn *et al.* 2020;  
389 Serrano-Solano *et al.*, 2022). For scientists, from a user's point of view, it  
390 offers extensive computing power and a graphical interface to use analysis  
391 workflows, even without experience in software development. Web-based  
392 access allows easy sharing of analytical workflows between collaborators and  
393 with a broader audience. Galaxy supports tools in almost any computational  
394 language, including R and Python, two of the most used languages in ecology,  
395 with many packages dedicated to ecological and biodiversity-oriented  
396 analyses incorporated (Lai *et al.*, 2019).

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

403 Galaxy is a powerful platform enabling researchers to readily move  
404 towards best practices. The Galaxy interface mitigates the difficulties  
405 associated with library management and code development, which permits  
406 simpler access to complex analytical methods. One can focus on the analysis  
407 itself and its concepts, rather than on syntax difficulties or cluster  
408 programming, disconnecting the study of data analysis concepts from the  
409 study of computing languages.

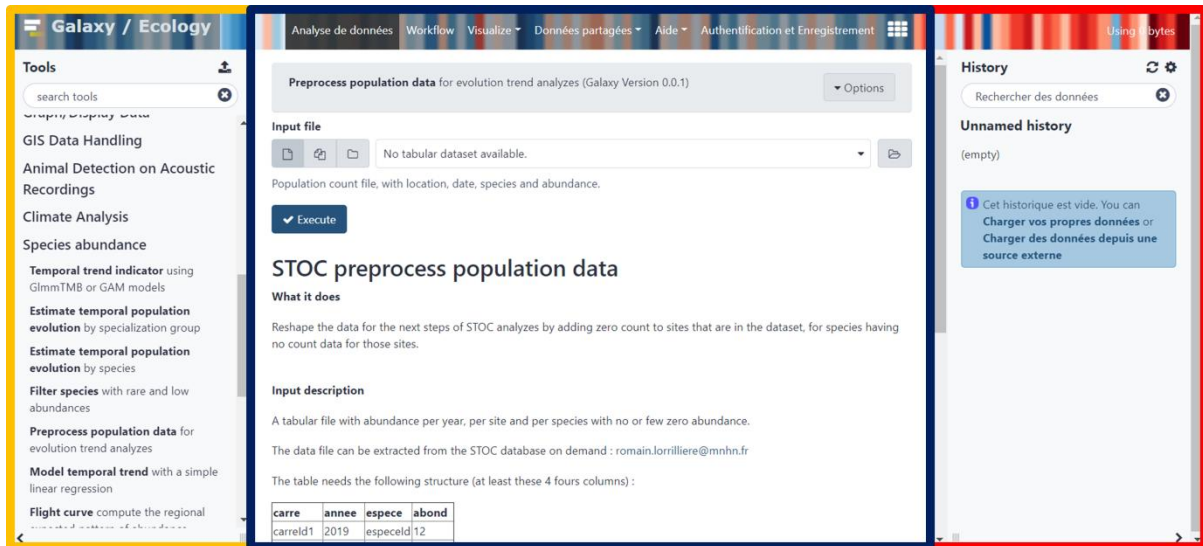
410 The platform is community-driven which permits continuous peer review of  
411 the platform and of the tools, workflows and tutorials provided. Many tutorials  
412 are available on the Galaxy Training Network (GTN) which is a valuable asset  
413 to the accessibility and reusability of tools and workflows (Batut *et al.*, 2018;  
414 Hiltmann *et al.*, 2023).

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

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

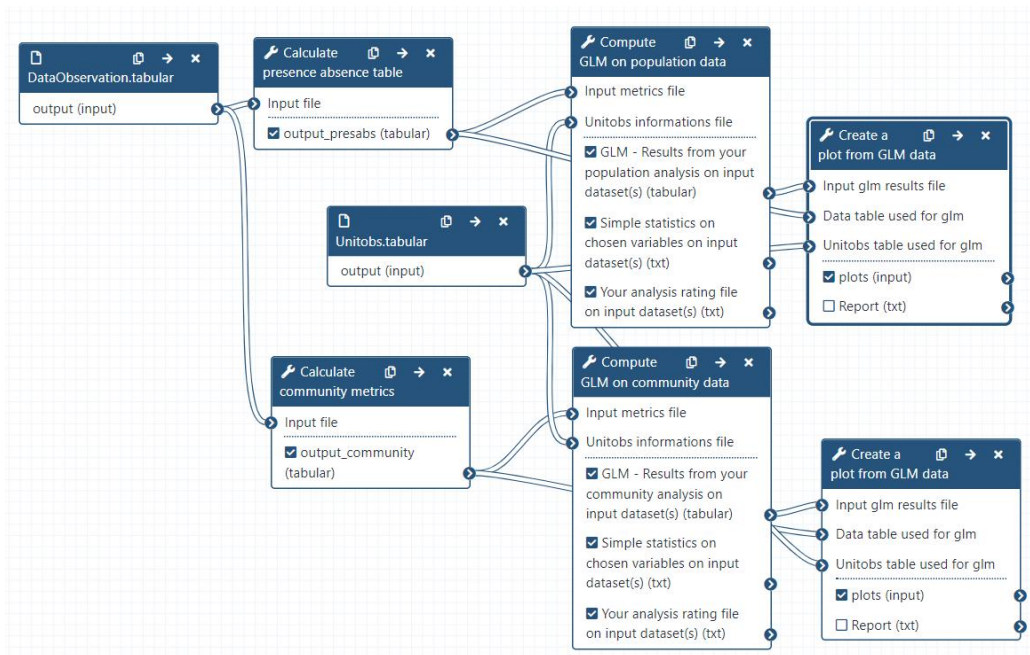
428  
429 Datasets can be uploaded on a Galaxy server from a local device, an  
430 online server, or a database. Users can then access every available tool (fig.

431 2, left panel) to modify, explore, and analyse their data. All tools used,  
 432 parameters, and data (inputs and outputs) of the analysis are saved in a  
 433 private “Galaxy history” (fig. 2, right panel), documenting every step of the  
 434 analytical procedure and recording the provenance of each output. From any  
 435 history, the user can extract a workflow (fig. 3) or directly share or publish  
 436 the history itself. Workflows are reusable through WorkflowHub  
 437 (<https://workflowhub.eu>) or Dockstore (<https://dockstore.org>) and exportable  
 438 in CWL and RO-CRATE standards.



439

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



443

444 **Figure 3** - Representation of a Galaxy workflow in the editing interface  
 445 of a Galaxy server. Each box represents an analysis tool, and the lines  
 446 represent the flow of data through the tools



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

452 Galaxy-Ecology has implemented workflows for biodiversity data  
453 exploration, eDNA processing, general population and community metrics  
454 and models, ecoregionalisation, NDVI (Normalised difference vegetation  
455 index) computation with Sentinel-2 data among others (see some examples:  
456 <https://workflowhub.eu/workflows/657>) and tutorials for several of them are  
457 available on the GTN platform (see [https://training.galaxyproject.org/training-](https://training.galaxyproject.org/training-material/topics/ecology)  
458 [material/topics/ecology](https://training.galaxyproject.org/training-material/topics/ecology)).

459 Eventually, one can modify or develop entirely new tools and workflows  
460 with any computational language to make them accessible to all users on any  
461 Galaxy server.

462 Galaxy is an utterly participative platform and several ways to participate  
463 to Galaxy exist depending on one's skills, available time, and needs. Anyone  
464 can participate to the Galaxy-Ecology initiative by notably:

- 465 – Sharing datasets, histories and workflows;
- 466 – Giving feedback on servers, tools, and workflows;
- 467 – Sharing tools and workflows ideas (eventually with code) through Git  
468 issues;
- 469 – Asking for tool modifications through issues;
- 470 – Modifying existing tools or proposing new tools through GitHub or  
471 GitLab;
- 472 – Writing or contributing to a GTN tutorial on a specific functionality or a  
473 workflow on the Galaxy Training Network platform;
- 474 – Create learning pathways, a set of tutorials curated by community  
475 experts to form a coherent set of lessons around a topic, building up  
476 knowledge ([https://training.galaxyproject.org/training-](https://training.galaxyproject.org/training-material/learning-pathways)  
477 [material/learning-pathways](https://training.galaxyproject.org/training-material/learning-pathways));
- 478 – Propose training events and help users in the utilisation of a workflow  
479 and tutorial.

480  
481 Analyses are rarely computed only once. Any analysis with a generalisation  
482 potential is a suitable candidate to be Galaxy-fied. A methodological  
483 framework is presented in online supplementary material  
484 ([https://github.com/ColineRoyaux/Galaxy\\_Templates/blob/main/Methods/Methods-](https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure_.md)  
485 [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 depending on  
486 potential interests, computing language skills, and willingness to invest more  
487 or less time in the process: (i) 'user' relying on existing Galaxy tools and  
488 workflows to analyse data (lower time investment), (ii) 'developer' relying on  
489 existing and validated analytical procedure to develop Galaxy tools and  
490 workflows (highest time investment), and (iii) 'trainer' relying on existing  
491 Galaxy tools to share workflows and create training material (variable time  
492 investment).  
493

495 As highlighted in previous sections, there are many best practices and  
496 recommendations existing for analytical procedures, data management, and  
497 computational code development. The levels of application of these best  
498 practices fall within a continuum offering many possibilities. From the lowest  
499 to the highest best practice levels for a published work there can be for  
500 example:

- 501 – Raw data and analytical procedure are not shared, only processed and  
502 interpreted results along with a brief description of methods.
- 503 – Pre-processed data is shared, and methods are described in the word-  
504 limit given by the publisher (example: tables of metrics and how it was  
505 calculated).
- 506 – Raw data and source code are shared on a repository. Software and  
507 package versions are not specified and there is no guaranty to be able  
508 to reproduce the analytical procedure.
- 509 – Raw data and atomised - generalised source codes are shared on a  
510 repository with specified hardware, software and dependencies  
511 versions. Input parameters are recorded in an attached file.
- 512 – Raw data is shared with proper metadata and an actionable version of  
513 the whole analytical procedure is traceable, ready to use and  
514 eventually reuse on other data types. Such level can be attained  
515 notably using Galaxy.
- 516 – All results and conclusions are published as an executable paper with  
517 analyses and workflows implemented and executable directly in the  
518 shared article (Strijkers *et al.*, 2011).

519 Executable Papers (Strijkers *et al.*, 2011) can require significant time and  
520 resource investment as well as good knowledge of programming languages,  
521 making it an admirable but hard-to-attain goal.

522 Atomisation and generalisation of computer codes can represent a  
523 relatively low investment strategy to attain certain levels of best practices  
524 such as transparency and reusability. It also carries advantages such as  
525 easier peer review, modularity of analytical procedures and, consequently,  
526 time savings. Indeed, applying the framework is not sufficient to attain the  
527 highest levels of best practices. For reproducibility and transparency, the  
528 management of the environment, softwares and package versions can be  
529 hard to maintain and record. A comprehensive tracking of input, outputs and  
530 codes requires meticulous management of files arborescence in the  
531 environment. Additionally, non-code developers will be able to partially  
532 review the analytical procedure only if the workflow is clearly outlined in an  
533 adapted format (e.g. table, graphical representation). Accessibility and  
534 findability of the atomised and generalised analytical procedure is dependent  
535 of its proper sharing (e.g. persistent link, open repository).

536 Galaxy can represent an easier gateway towards higher levels of best  
537 practice as sharing a complete, detailed and (re-)executable analytical  
538 procedure is facilitated through provenance tracking and automatic metadata  
539 enrichment. In comparison, many scientific workflow management systems,  
540 such as Snakemake, Nextflow or the R package Targets, operate from the



541 command line. In ecology, numerous initiatives have tried to introduce such  
542 systems, starting with more user-friendly solutions. For example, the KNIME  
543 and Kepler systems with the CoESRA initiative (Collaborative Environment for  
544 Scholarly Research and Analysis) in Australia; Taverna with the BioVeL  
545 initiative (Biodiversity Virtual e-Laboratory) in Europe; or very recently, the  
546 BON in a Box pipeline engine. These systems are more accessible to new  
547 users by offering a graphical interface while achieving high specificity  
548 (Berthold *et al.*, 2007; Hardisty *et al.*, 2016; <https://boninabox.geobon.org/>).  
549 However, good computer programming or scientific workflow management  
550 knowledge is still necessary to use these applications correctly.

551 In comparison to the atomisation-generalisation framework, Galaxy can be  
552 rightfully seen as heavier for experienced programmers as it requires to learn  
553 to use a new platform. Additionally, more effort may be required on Galaxy  
554 when an additional analytical step needs to be developed, but the Galaxy  
555 community can be an efficient crutch on which hard-pressed scientists can  
556 rely. Indeed, one can ask for help on the implementation of tools whether one  
557 knows computing languages and can share their code or not.

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

## 571 Acknowledgements

572 Authors want to thank Sandrine Pavoine for its highly relevant and helpful  
573 advices and reviews on both the content and the form of the article.

### 574 Authors contribution statement

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

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

578 J.-B. M. and Y. L.B. reviewed and helped rewrite many parts of the draft.

579 Y. R. and D. P. helped inspire and were invested in the early design of the  
580 article.

581 M. J. and P. S. tested and approved the appliance of the framework.

582 O. N., M. J., Y. R., M. E., B. B., A. F., H. R. and S. H. highly enhanced the  
583 quality of the redaction in both form and content at several stages of the  
584 draft.

585 H. R., S. H., B. B., A. F., and B. G. are involved in the Galaxy-E initiative and  
586 provided many advices on the redaction of the article and/or on the  
587 development of the initiative.

588 M. E. and G. M. are involved in Antarctic-oriented Galaxy tool and workflow  
589 development coordination.

590 C. B., R. L., A. M., Y. B., A. A., T. V. and V. C. developed scripts, tools  
591 and/or Galaxy workflows to contribute to the Galaxy-E initiative.

592 E. A. developed R scripts and apps used to integrate R Shiny apps as  
593 Galaxy interactive tools and initiate "Research Data management Galaxy  
594 tools".

595 E. M. and C. U. developed the first training materials for Galaxy-E.

596 E. T. worked on the use of the first Galaxy-E analysis.

597 M. D., G. L. and R. J. were coordinating the prefiguration of Galaxy-E  
598 through the 65 Millions d'Observateurs project.

599 Additionnally, all authors reviewed and approved the article draft.

600

## Funding

601 Funding were provided by the European Union through the Erasmus+  
602 project; the Agence Nationale de la Recherche through the 65 Million  
603 d'Observateurs and the IA-Biodiv projects; the French National Fund for Open  
604 Science through the OpenMetaPaper project; the European commission  
605 through the H2020, the EOSC-Pillar, and the H2020 GAPARS projects; the GO  
606 FAIR initiative through the BiodiFAIRse Implementation Network; the Blue  
607 Nature Alliance; and the Antarctic and Southern Ocean Coalition. Finally,  
608 funding by the French Ministry of Higher Education and Research were  
609 provided for the "Pôle national de données de biodiversité" e-infrastructure.

610

## Conflict of interest disclosure

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

613

## References

- 614 Araújo MB, Anderson RP, Barbosa AM, Beale CM, Dormann CF, Early R, Garcia  
615 RA, Guisan A, Maiorano L, Naimi B, O'Hara RB, Zimmermann NE, Rahbek C  
616 (2019) Standards for distribution models in biodiversity assessments.  
617 *Science Advances*, **5**, 1-12. <https://doi.org/10.1126/sciadv.aat4858>
- 618 Archmiller AA, Johnson AD, Nolan J, Edwards M, Elliott LH, Ferguson JM,  
619 Iannarilli F, Vélez J, Vitense K, Johnson DH, Fieberg J (2020) Computational  
620 Reproducibility in The Wildlife Society's Flagship Journals. *Journal of*  
621 *Wildlife Management*, **84**, 1012-1017. <https://doi.org/10.1002/JWMG.21855>
- 622 Batut B, Hiltemann S, Bagnacani A, Baker D, Bhardwaj V, Blank C,  
623 Bretaudeau A, Brillet-Guéguen L, Čech M, Chilton J, Clements D, Doppelt-  
624 Azeroual O, Erxleben A, Freeberg MA, Gladman S, Hoogstrate Y, Hotz HR,  
625 Houwaart T, Jagtap P, Larivière D, Le Corguillé G, Manke T, Mareuil F,  
626 Ramírez F, Ryan D, Sigloch FC, Soranzo N, Wolff J, Videm P, Wolfien M,  
627 Wubuli A, Yusuf D, Taylor J, Backofen R, Nekrutenko A, Grüning B (2018)

628 Community-Driven Data Analysis Training for Biology. *Cell Systems*, **6**,  
629 752-758. <https://doi.org/10.1016/j.cels.2018.05.012>

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

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

636 Boyd RJ, August TA, Cooke R, Logie M, Mancini F, Powney GD, Roy DB, Turvey  
637 K, Isaac NJB (2023) An operational workflow for producing periodic  
638 estimates of species occupancy at national scales. *Biological Reviews*, **98**,  
639 1492–1508. <https://doi.org/10.1111/brv.12961>

640 Carroll S, Garba I, Figueroa-Rodríguez O, Holbrook J, Lovett R, Materechera S,  
641 Parsons M, Raseroka K, Rodriguez-Lonebear D, Rowe R, Sara R, Walker J,  
642 Anderson J, Hudson M (2020) The CARE Principles for Indigenous Data  
643 Governance. *Data Science Journal*, **19**, 43. [https://doi.org/10.5334/dsj-](https://doi.org/10.5334/dsj-2020-043)  
644 [2020-043](https://doi.org/10.5334/dsj-2020-043)

645 Casajus N. (2023) {rcompendium} {An} {R} package to create a package or  
646 research compendium structure.

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

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

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

660 Di Cosmo R, Zacchiroli S (2017) Software Heritage: Why and How to Preserve  
661 Software Source Code.

662 Di Tommaso P, Chatzou M, Floden EW, Barja P., Palumbo E, Notredame C  
663 (2017) Nextflow enables reproducible computational workflows. *Nature*  
664 *Biotechnology*, **35**, 316–319. <https://doi.org/10.1038/nbt.3820>

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

667 EMBL Australia Bioinformatics Resource (2013) Community Survey Report  
668 [https://www.embl-abr.org.au/news/braembl-community-survey-report-](https://www.embl-abr.org.au/news/braembl-community-survey-report-2013/)  
669 [2013/](https://www.embl-abr.org.au/news/braembl-community-survey-report-2013/)

670 Emery NC, Crispo E, Supp SR, Farrell KJ, Kerkhoff AJ, Bledsoe EK, O'Donnell KL,  
671 McCall AC, Aiello-Lammens ME (2021) Data Science in Undergraduate Life  
672 Science Education: A Need for Instructor Skills Training. *BioScience*, **71**,  
673 1274–1287. <https://doi.org/10.1093/BIOSCI/BIAB107>

674 European Commission, Directorate-General for Research and Innovation  
675 (2018) Cost-benefit analysis for FAIR research data: cost of not having  
676 FAIR research data. *Publications Office*. <https://doi.org/10.2777/02999>  
677 Fanelli D (2018) Is science really facing a reproducibility crisis, and do we  
678 need it to? *Proceedings of the National Academy of Sciences of the United*  
679 *States of America*, **115**, 2628–2631.  
680 <https://doi.org/10.1073/pnas.1708272114>  
681 Fang FC, Casadevall A (2015) Competitive Science: Is Competition Ruining  
682 Science? *Infection and Immunity*, **83**, 1229–1233.  
683 <https://doi.org/10.1128/IAI.02939-14>  
684 Farley SS, Dawson A, Goring SJ, Williams JW (2018) Situating Ecology as a  
685 Big-Data Science: Current Advances, Challenges, and Solutions. *BioScience*,  
686 **68**, 563–576. <https://doi.org/10.1093/BIOSCI/BIY068>  
687 Field B, Booth A, Illott I, Gerrish K (2014) *Using the Knowledge to Action*  
688 *Framework in practice: a citation analysis and systematic review.*  
689 *Implementation Science*, **9**, 172. [https://doi.org/10.1186/s13012-014-0172-](https://doi.org/10.1186/s13012-014-0172-2)  
690 [2](https://doi.org/10.1186/s13012-014-0172-2)  
691 Goble C, Cohen-Boulakia S, Soiland-Reyes S, Garijo D, Gil Y, Crusoe MR,  
692 Peters K, Schober D (2020) FAIR Computational Workflows. *Data*  
693 *Intelligence*, **2**, 108–121. [https://doi.org/10.1162/dint\\_a\\_00033](https://doi.org/10.1162/dint_a_00033)  
694 Gomes DGE, Pottier P, Crystal-Ornelas R, Hudgins EJ, Foroughirad V, Sánchez-  
695 Reyes LL, Turba R, Martinez PA, Moreau D, Bertram MG, Smout CA, Gaynor  
696 KM (2022) Why don't we share data and code? Perceived barriers and  
697 benefits to public archiving practices. *Proceedings of the Royal Society B*,  
698 **289**, 20221113 <https://doi.org/10.1098/rspb.2022.1113>  
699 Gownaris NJ, Vermeir K, Bittner MI, Gunawardena L, Kaur-Ghumaan S,  
700 Lepenies R, Ntsefong GN, Zakari IS (2022) Barriers to Full Participation in  
701 the Open Science Life Cycle among Early Career Researchers. *Data*  
702 *Science Journal*, **21**, 2. <https://doi.org/10.5334/DSJ-2022-002>  
703 Grüning B, Chilton J, Köster J, Dale R, Soranzo N, van den Beek M, Goecks J,  
704 Backofen R, Nekrutenko A, Taylor J (2018) Practical Computational  
705 Reproducibility in the Life Sciences. *Cell Systems*, **6**, 631–635.  
706 <https://doi.org/10.1016/j.cels.2018.03.014>  
707 Hampton SE, Jones MB, Wasser LA, Schildhauer MP, Supp SR, Brun J,  
708 Hernandez RR, Boettiger C, Collins SL, Gross LJ, Fernández DS, Budden A,  
709 White EP, Teal TK, Labou SG, Aukema JE (2017) Skills and Knowledge for  
710 Data-Intensive Environmental Research. *BioScience*, **67**, 546–557.  
711 <https://doi.org/10.1093/BIOSCI/BIX025>  
712 Hardisty AR, Bacall F, Beard N, Balcázar-Vargas MP, Balech B, Barcza Z,  
713 Bournalat SJ, Giovanni R, Jong Y, Leo F, Dobor L, Donvito G, Fellows D, Guerra  
714 AF, Ferreira N, Fetyukova Y, Fosso B, Giddy J, Goble C, Güntsch A, Haines R,  
715 Ernst VH, Hettling H, Hidy D, Horváth F, Ittzés D, Ittzés P, Jones A,  
716 Kottmann R, Kulawik R, Leidenberger S, Lyytikäinen-Saarenmaa P, Mathew  
717 C, Morrison N, Nenadic A, Hidalgo AN, Obst M, Oostermeijer G, Paymal E,  
718 Pesole G, Pinto S, Poigné A, Fernandez FQ, Santamaria M, Saarenmaa H,  
719 Sipos G, Sylla KH, Tähtinen M, Vicario S, Vos RA, Williams AR, Yilmaz P  
720 (2016) BioVeL: A virtual laboratory for data analysis and modelling in

721 biodiversity science and ecology. *BMC Ecology*, **16**, 49.  
722 <https://doi.org/10.1186/S12898-016-0103-Y>

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

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

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

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

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

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

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

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

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

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

767 Larcombe L, Hendricusdottir R, Attwood T, Bacall F, Beard N, Bellis L, Dunn W,  
768 Hancock J, Nenadic A, Orengo C, Overduin B, Sansone S, Thurston M, Viant



769 M, Winder C, Goble C, Ponting C, Rustici G (2017) ELIXIR-UK role in  
770 bioinformatics training at the national level and across ELIXIR.  
771 *F1000Research*, **6**, 952. <https://doi.org/10.12688/f1000research.11837.1>  
772 Leroy B (2023) Choosing presence-only species distribution models. *Journal of*  
773 *Biogeography*, **50**, 247–250. <https://doi.org/10.1111/jbi.14505>  
774 Lortie CJ (2021) The early bird gets the return: The benefits of publishing your  
775 data sooner. *Ecology and Evolution*, **11**, 10736–10740.  
776 <https://doi.org/10.1002/ECE3.7853>  
777 McIntire EJB, Chubaty AM, Cumming SG, Andison D, Barros C, Boisvenue C,  
778 Haché S, Luo Y, Micheletti T, Stewart FEC (2022) PERFICT: A Re-imagined  
779 foundation for predictive ecology. *Ecology Letters*, **25**, 1345–1351.  
780 <https://doi.org/10.1111/ELE.13994>  
781 Michener WK (2015) Ten Simple Rules for Creating a Good Data Management  
782 Plan. *PLOS Computational Biology*, **11**, e1004525.  
783 <https://doi.org/10.1371/JOURNAL.PCBI.1004525>  
784 Michener WK, Jones MB (2012) Ecoinformatics: Supporting ecology as a data-  
785 intensive science. *Trends in Ecology and Evolution*, **27**, 85–93.  
786 <https://doi.org/10.1016/j.tree.2011.11.016>  
787 Minocher R, Atmaca S, Bavero C, McElreath R, Beheim B (2021) Estimating  
788 the reproducibility of social learning research published between 1955 and  
789 2018. *Royal Society Open Science*, **8**, 210450.  
790 <https://doi.org/10.1098/RSOS.210450>  
791 Munafò MR, Nosek BA, Bishop DVM, Button KS, Chambers CD, Percie Du Sert  
792 N, Simonsohn U, Wagenmakers EJ, Ware JJ, Ioannidis JPA (2017) A  
793 manifesto for reproducible science. *Nature Human Behaviour*, **1**, 0021.  
794 <https://doi.org/10.1038/s41562-016-0021>  
795 Natural Environment Research Council (2010, 2012) Most Wanted:  
796 Postgraduate Skills Needs in the Environment Sector.  
797 Plesser HE (2018) Reproducibility vs. Replicability: A brief history of a  
798 confused terminology. *Frontiers in Neuroinformatics*, **11**, 76.  
799 <https://doi.org/10.3389/FNINF.2017.00076>  
800 Powers SM, Hampton SE (2019) Open science, reproducibility, and  
801 transparency in ecology. *Ecological applications*, **29**, e01822.  
802 <https://doi.org/10.1002/eap.1822>  
803 Serrano-Solano B, Fouilloux A, Eguinoa I, Kalaš M, Grüning B, Coppens F  
804 (2022) Galaxy: A Decade of Realising CWFR Concepts. *Data Intelligence*, **4**,  
805 358–371. [https://doi.org/10.1162/dint\\_a\\_00136](https://doi.org/10.1162/dint_a_00136)  
806 Soiland-Reyes S, Sefton P, Crosas M, Castro LJ, Coppens F, Fernández JM,  
807 Garijo D, Grüning B, La Rosa M, Leo S, Ó Carragáin E, Portier M, Trisovic A,  
808 Community R-C, Groth P, Goble C (2022) Packaging research artefacts with  
809 RO-Crate. *Data Science*, **5**, 97–138. <https://doi.org/10.3233/DS-210053>  
810 Strijkers R, Cushing R, Vasyunin D, De Laat C, Belloum ASZ, Meijer R (2011)  
811 Toward executable scientific publications. *Procedia Computer Science*, **4**,  
812 707–715. <https://doi.org/10.1016/j.PROCS.2011.04.074>  
813 The Galaxy Community (2022) The Galaxy platform for accessible,  
814 reproducible and collaborative biomedical analyses: 2022 update. *Nucleic*  
815 *acids research*, **50**, W345–W351. <https://doi.org/10.1093/NAR/GKAC247>

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

819 Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A,  
820 Blomberg N, Boiten JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes  
821 AJ, Clark T, Crosas M, Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R,  
822 Gonzalez-Beltran A, Gray AJG, Groth P, Goble C, Grethe JS, Heringa J, t  
823 Hoen PAC, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ, Martone ME, Mons A,  
824 Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R, Sansone SA,  
825 Schultes E, Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, Van  
826 Der Lei J, Van Mulligen E, Velterop J, Waagmeester A, Wittenburg P,  
827 Wolstencroft K, Zhao J, Mons B (2016) Comment: The FAIR Guiding  
828 Principles for scientific data management and stewardship. *Scientific Data*,  
829 **3**, 1-9. <https://doi.org/10.1038/sdata.2016.18>

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

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