

Authors want to deeply thank the recommender as well as the two reviewers for their suggestions and advices. We believe the manuscript has been significantly enhanced by their perspectives.

Decision for round #1 : *Revision needed*

Revision needed to improve readability

Both reviewers see the value of the work (and I agree), but make important comments that would improve the readability of the manuscript. I think the separation between philosophy, implementation/architecture, and anticipated use-cases needs to be clarified. This will help readers navigate a very dense manuscript, especially readers with different types of background/expertise.

by *Timothée Poisot*, 07 Jun 2024 12:38

Manuscript: <https://doi.org/10.32942/X2G033>

version: 1

Thank you, we reorganised the manuscript to improve the readability. Notably, as the reviewers seemed to appreciate the visual abstract, we followed the same organisation for the manuscript and referenced it with numbers in the body of the manuscript (lines 171, 247-248, 289, 368). Additionally, many superfluous paragraphs were removed and in each section of the manuscript we separated paragraphs focused on concepts and philosophy, and paragraphs on implementation. The presented framework and the Galaxy-E platform can host many use-cases, each have been developed to be applicable to a large variety of analytical procedures. Some examples of workflows developed on Galaxy are presented at lines 452-458, as these examples are extensively described at the indicated links, we did not describe them deeply in the manuscript to avoid further densification of the manuscript. However, we are open to make a deeper description if you believe it will help to enhance comprehension. Finally, we transferred the methods section as online supplementary material to update it if needed

(https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure.md).

Please, do not hesitate if other alterations are needed for the clarity of this manuscript, we are eager to make sure it is accessible to the largest public possible.

Review by anonymous reviewer 1, 23 May 2024 15:29

This is a paper in three parts. The first part discusses issues around analytical practices in ecology, and the principles by which these can become more reproducible. The second half is a description of the Galaxy Ecology platform, and it's potential to realise better practice among the community. The third part described the technical details of implementation.

The visual abstract is really great, presenting a clear, cohesive message about how atomisation and generalisation can improve reproducibility and FAIR principles. However, the main text does not reflect this clarity of message.

Thank you, we reorganised the manuscript in accordance with the visual abstract, and referenced each section of the manuscript with the corresponding part of the visual abstract to enhance clarity and links between the manuscript and visual abstract (lines 171, 247-248, 289, 368).

Overall there is a lot of good material in here, but the article does not feel mature in how it is structured and presented. I do not get a clear idea who the article is aimed at, nor who is the

target user of Galaxy Ecology. It's hard to see how the present version would lead to material change in how ecologists go about their work. Another challenge is that much of the material about Galaxy Ecology is likely to become quickly outdated. For these reasons, my overall suggestion would be to greatly reduce the amount of text in the section on Galaxy and in the Methods. Rewrite these sections in a way that uses Galaxy to illustrate the general principles in the first part of the manuscript.

This article is aimed at researchers and experts in ecology that use analytical tools (R or else). The target user of Galaxy can be many publics depending on their interests and skills. We added a paragraph to describe possible users in the Galaxy-E section (lines 369-375).

We reduced, reorganised and rewrote all sections about Galaxy.

The paper needs to make more comparisons to current workflow platforms. The examples provided are not current: Taverna retired in 2020 and the latest version of Kepler (version 2.5) was 2015. For example, ecologically specific tools such as Bon-in-a-box <https://boninabox.geobon.org/> or more general tools such as the R package Targets: <https://books.ropensci.org/targets/>.

Thank you, we didn't think of comparing to BON in a Box as it has a wider purpose (networking, reporting, identifying data gaps) and the pipeline engine has been released very recently but both suggestions are now included (lines 539-550). After looking into tutorials for using these two solutions, we however get to same conclusions as for other examples we previously cited, Targets requires knowing R language and the installation and use of BiaB is (for now) requiring to use command line, git and Docker which is not as easy to use as an online platform.

Additionally, the methods section could be better as online supplementary materials (or zenodo repository) of simply part of the galaxy user guide. It is not relevant to the core message of the manuscript. As the galaxy-E platform develops, the user guide may change. Hosting this content online rather than in the manuscript allows that updating to ensure parity with the Galaxy-E platform.

Thank you, it is a good idea. This section has been removed and added as a markdown file here:

https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure_.md

The introduction needs to be more concise, a lot of text is used to set up the wider concerns about reproducibility in ecology, which is important, but the main contribution is not to review ecological reproducibility to provide the solution and at current it takes too long to get to the solution. Large parts of the text are not necessary for the delivery of the core message eg. lines 106-131, 151-185, 211-229. In the detailed comments below, I've made a few suggestions for how the Introductory sections could help to contextualise the issues being discussed.

Thank you for your suggestions, we reduced the Introduction by ~20%. For lines 106-131 (107-113 in the new version) and lines 151-185 (138-155 in the new version), we removed most sentences apart from the definitions of reproducibility, workflow and computational reproducibility.

For lines 211-229 (181-194 in the new version), we removed approximately half of the paragraph except from the explanation of why presented best practices are important for science in general along with the transition on the growing complexity of analytical procedures. This part permits to demonstrate the need for a simplified representation of analyses.

For the section on Galaxy, I think the authors need to give a much clearer exposition of what it is, who it is for and how it can help to deliver the principles outlined in the first section. At present, this text assumes too much knowledge of the system for a naïve reader to properly engage.

We reorganised the section and added a few sentences to clearly state the aim and scope of Galaxy (lines 369-375). We have structured the section to include paragraphs on the philosophy of Galaxy at the beginning, followed by those on implementation and architecture.

The Methods section is even more difficult to follow: it's half-way between a user-manual and a conceptual overview but doesn't quite deliver to either of those goals. Perhaps it would help if these issues were illustrated via the use of one or more case studies.

This section has been removed and added as a markdown file here: https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure_.md.

Detailed comments:

· Abstract Line 75: perhaps add a few words to explain that the principles described here are applicable across all levels of ecological analysis, e.g. “from individual research projects to production-level analytical pipelines”

Thank you, done (line 76 in the new version).

· Abstract line 77: Perhaps explain why “atomisation” is the right word. I thought this was a typo and that the authors meant to write “automation” instead. The authors need to provide rationale behind using the term ‘atomisation’ to describe the process of making a script more modular, there is no previous usage of this term in the software development literature. Whereas modular programming is a widely used and understood term (https://en.wikipedia.org/wiki/Modular_programming). If the authors want to use the term ‘atomisation’ it must be made clear it is a new, non-standard term introduced in this manuscript that may not be widely understood across other sectors.

Yes, it is true that “atomisation” we present here can be linked to modular programming in their values and goals. However, modular programming is a software development technique dedicated to computer scientists and can have various granulometries. Depending on languages and aims, it can refer to a function, or to a library, or to a package etc... Here, atomisation is dedicated to a narrower subject that is analytical procedures but has a wider meaning. It refers to the separation of an analytical workflow in single tasks/steps, it is not necessarily programming-related. As it is presented in the PERFICT article we cite (McIntire et al., 2022), ecologists use modules, functions, packages and libraries to build their analytical workflows and the “atom” we present is the task they are willing to achieve with these modules, functions, packages and libraries. It is true that “modularisation” can be used but we chose a different term to avoid confusion.

We added clearer statement of the definition of atomisation and generalisation in the abstract (lines 77-82) and the Introduction (lines 219-224).

· Introduction paragraph 2: a useful citation here would be Cassey & Blackburn, who distinguished between “reproducibility” and “repeatability”, and discussed the relative merits of each. Also, with reference to later discussion of reproducibility around lines 205, it would be good to acknowledge that reproducibility is not an absolute but rather a relative concept (who needs the results to be identical to the 9th significant figure?).

Thank you, it is indeed important to state that reproducibility has been a long confused term. As this article is oriented towards analysis, we decided to use the article of Cohen-Boulakia et al. (2017) that gives distinctions between “repeatability”, “replicability”, “reproducibility” and “reusability” in the context of computational reproducibility. Here, reproducibility is not about getting the exact same results but more about getting results leading to similar conclusions (lines 107-113). We shortened the paragraph and clarified these issues.

· Lines 117-119: complicated sentence. I suggest to simplify to “Given the increasing complexity of ecological analyses, there is a clear need for tools that facilitate greater reproducibility.”

The sentence has been deleted.

· The same paragraph (lines 119-130, but also the next section, starting at line 140) would really benefit from a big more context for what is the problem and why atomised workflows are needed. My perspective is that ecology has, until now, been a discipline in which most analysis happened on a single computer, but increasingly we are seeing papers derived from big collaborations involving code developed in different labs. This means we are moving into an era where analytical pipelines are becoming so complex that no individual researcher can understand all the details at a granular level. Other disciplines (e.g. meteorology, particle physics) have already passed through this phase. There is plenty of literature that could help illustrate these issues and clarify why this paper is novel. First, it would be good to include a citation to support the assertion that analyses are becoming more complex. One option would be to find some data on the average length of supplementary information on ecological papers (I know of no such data!). An alternative would be to cite papers that describe highly complex workflows, e.g. Boyd et al (2023) or Jetz et al (2019).

Thank you, we also added the article from Leroy (2022) on the decision process behind presence-only SDMs. As most of the paragraph commented here has been removed and this particular matter is already partly discussed at lines 227-233 (188-198 in the new version), we transferred this suggestion into these lines.

Second, it might be worth acknowledging that individual branches of ecology have developed principles to enhance reproducibility within those sub-domains. The SDM community is perhaps the best example: citing papers such as Araujo et al (2019) and Golding et al (2017) might provide a way for the authors to explain what is different about the proposals in this paper.

We also added the article from Zurell et al. (2020) on the ODMAP protocol to report SDMs (lines 134-137). Unfortunately, the zoon package described by Golding et al. (2017) is not maintained anymore so we did not include this citation.

· Line 156: try to avoid directly quoting from another article

Quoting removed (lines 138-141).

· line 170: I would question whether long-term public archiving of code is as valuable as the authors assert. The most popular coding language among ecologists, R, is in a continuous state of evolution. Most R code written 10 years ago would not execute today. I’m not saying that we should not archive code: I just think it is important to be clear about what we are trying to achieve as a community and make decision about where to invest resources accordingly.

Indeed, long-term public archiving is necessary but insufficient to achieve computational reproducibility. It is not guaranteed to re-execute a R code 10 years after its development but facilitated if described through agnostic code recipe as conda and encapsulated to a dedicated container image. However, “time proofing” the execution is not the main objective of code archiving, it is mostly to keep a detailed track of what type of analyses has been used to produce a given result that lead to a given conclusion.

Nevertheless, it is true that this statement is not directly serving the message of the article, we decided to remove the mention of long term archiving in favour of the wider term “sharing” (line 142).

· Line 195: in the previous paragraph you made the case that code should be considered as data. So, for clarity, insert the word “observational” before “data” in this sentence.

Thank you, done (line 165).

· Line 298: “Atomisation refers to dividing ...”

Thank you, done (line 246).

- Line 298-366: this text on atomisation and generalisation is absolutely fine, but I can't help thinking that these must be fundamental principle of computer science. If so, it may be worth mentioning here, perhaps with a citation.

As you rightfully proposed earlier, it is true that these concepts can be linked to existing computer science principles, however it is not exactly the same idea. We do not want to induce confusion between computer science oriented principles and the atomisation-generalisation framework that is more domain science oriented. We absolutely do not intend to attribute ourselves the "invention" of a fundamental principle, that is why we cited the work of McIntire et al. (2022) that is ecology-oriented (line 257-270). We agree that it is often unnecessary to "create" new terminologies but it also seems to us that adding a new definition to an already multi-definition technique such as "modular programming" and "modules" can add vagueness to the point. For generalisation, we were not able to find if it is in fact a fundamental principle, it is discussed in some online posts and tutorials (<https://medium.com/the-art-of-software-development/generalization-a-key-technique-in-programming-c0e71166d98e>) with a similar definition as we use in this article but we found nothing in the academic literature. As we are mostly ecologists and not computer scientists we are a bit hesitant to make comparisons with terms from a domain we do not master as we might misinterpret them.

If it is still not clear enough in the text of the article, please do not hesitate to provide us with any additional suggestion.

- Line 369: "et" -> "at"

Thank you, done (line 304).

- Line 386: missing word "A" at beginning of sentence

Thank you, done (line 321).

- Line 412: at the beginning of this section, it would be useful to explain who is the target user. Are you recommending that everyone in ecology use it for all of their analysis? Or is it better suited to large collaborative projects?

We added a paragraph to describe possible users in this section (lines 369-375). Galaxy is suited for individual analyses as well as for large collaborative projects. However, we do not necessarily recommend or think everyone in ecology should use Galaxy for their analyses because we cannot possibly have a holistic vision on what is suited for everyone. Here, we are briefly presenting the Galaxy-E solution and if it is relevant to enhance the best practice application for some cases we are satisfied.

- Line 435: "tools" -> "tool"

Thank you, done (line 429).

REFERENCES

Araujo et al (2019) Standards for distribution models in biodiversity assessments. Science Advances <https://www.science.org/doi/full/10.1126/sciadv.aat4858>

Boyd et al (2023) Biological Reviews <https://onlinelibrary.wiley.com/doi/full/10.1111/brv.12961>

Cassey & Blackburn (2006) Reproducibility and Repeatability in Ecology. BioScience <https://academic.oup.com/bioscience/article/56/12/958/221622>

Golding et al (2017) the Zoon package ... Methods in ecology & Evolution. <http://doi.wiley.com/10.1111/2041-210X.12858>

Jetz et al (2019) Nature Ecology & Evolution <https://doi.org/10.1038/s41559-019-0826-1>

Cohen-Boulakia S et al. (2017) Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities. *Future Generation Computer Systems*. <https://doi.org/10.1016/j.future.2017.01.012>

Leroy B (2023) Choosing presence-only species distribution models. *Journal of Biogeography*. <https://doi.org/10.1111/jbi.14505>

McIntire EJB et al. (2022) PERFICT: A Re-imagined foundation for predictive ecology. *Ecology Letters*. <https://doi.org/10.1111/ELE.13994>

Zurell D et al. (2020) A standard protocol for reporting species distribution models. *Ecography*. <https://doi.org/10.1111/ecog.04960>

Review by anonymous reviewer 2, 09 May 2024 19:15

This manuscript introduces and describes the Galaxy-Ecology tool, laying out different modes of engagement and the ways the approach addresses reproducibility issues in ecology. The Galaxy-Ecology project and associated community seem like a powerful framework to build and share ecological analyses and this paper includes all the essential introductory pieces for a new user. I particularly appreciate the way the paper discusses different types of users and the value it adds for each. I do however think there is a fair bit of related but ultimately superfluous content that, if paired down, would greatly improve the readability and clarity of the paper.

While the initial discussion of reproducibility is clearly motivating for development of the tool, the current level of depth is unnecessary and even a little misleading for the reader. For example, Galaxy-Ecology isn't mentioned until the fourth section of the introduction and is introduced in a way that makes the reader unsure if it's just a nice example to illustrate the reproducibility point or the main message of the paper in and of itself. Given that a history of reproducibility in ecology is not the goal of the paper, I would recommend editing everything prior to "Framework towards good practices" down to a few introductory paragraphs that also immediately introduce Galaxy-Ecology as a solution. If I have misunderstood the purpose of the paper and the goal of the paper is to first give a detailed context for the reproducibility crisis and second discuss Galaxy-Ecology, that structure can be better set up in the abstract and first couple paragraphs of the paper.

We altered the paragraphs in the Introduction about reproducibility and best practices in general to only keep contextualisation and definition sentences that serve the key messages of the manuscript.

Here, the purpose is to propose several actionable ways to achieve better best practice. The framework atomisation-generalisation is the simplest and lowest time investment which is, in our opinion, the most relevant to ecology researchers as most don't have much time to learn to use or to contribute to Galaxy-E. Galaxy-E is also a simple solution and permits to attain higher levels of best practice but can in some cases (e.g. the needed workflow is not available on the platform and needs to be developed) require a lot of investment. We are planning to propose another article that is entirely Galaxy-oriented which is why it is, in this article, more used to illustrate which best practices can be attained and how. We added elements about the framework and reduced paragraphs about Galaxy in the Introduction to clarify the aim of the article (lines 212-243).

In general I found the structure of the paper hard to follow, as technical details about engaging with the tool or development are interspersed with motivation and philosophy. One possible approach for addressing that confusion is to lay out the technical details in an initial description of the tool, then describe which of those pieces different kinds of users might engage with, rather than introducing new information in the user sections.

The structure of the article has been modified in accordance with the visual abstract, and referenced each section of the manuscript with the corresponding part of the visual abstract to enhance clarity and links between the manuscript and

visual abstract (lines 171, 247-248, 289, 368). In the section about Galaxy, we separated paragraphs focused on concepts and philosophy (lines 362-426), and paragraphs on implementation (lines 428-493).

As the reviewer #1 suggested, the methods section is susceptible to evolve in the future so we changed this section into online supplementary material (https://github.com/ColineRoyaux/Galaxy_Templates/blob/main/Methods/Methods%20-%20How%20to%20Galaxy-fy%20your%20analytical%20procedure.md) that is possible to update. However, on your advice, we moved some elements from the methods section to the Galaxy-E section.

In a related formatting issue, currently the three guidelines sections read as multi-paragraph lists which I found fairly jarring. Rather than referencing different steps of the workflow as a starting bullet point, I would reference steps in the body of the paragraph (for example as “(step A)”) to aid flow of the writing.

At first, the methods section was written as a continuous paragraph and steps were referred as you suggested in the text. Unfortunately, the separation between each step was very unclear in this format and, as we tested the method, it seemed hard to follow the described procedures in this format. We tried to change this section back to this format as you suggested and it was not didactic as it should be. We tried several formats but were not able to find a better compromise. We hope, as it is now an online supplementary material, that it will not be a problem for you.

As a small linguistic note, I would suggest the phrase “good practices” be replaced with “best practices”, which is a more common way to reference standards and will be immediately recognizable to a reader.

Thank you, done in the whole manuscript.

A few line comments:

Line 461: Unnecessary reference, either remove sentence or integrate into a paragraph.

Reference removed.

Line 560: This heading is unnecessary.

Header removed.

Line 632: The colon here is confusing to me, I can see from the workflow what you're implying but the list within a list is quite hard to follow in the writing. This would be much easier to express in a paragraph rather than list format.

Thank you, the modification has been included in the online supplementary material.

Line 644: This is a good example of the kind of concept that should be introduced outside the user descriptions. Is the "Galaxy history" an internal versioning system? At what level is it tracking, just the kinds of modules that are used in the pipeline?

The Galaxy history is described on lines 432-433 and in Table 2. Some advantages of the Galaxy history are described afterwards. The user can open as many histories as needed for example one history opened for an analysis at time T and another one for an update of the same analysis at time $T+1$. As the histories are not directly linked to each other it is not technically a versioning system but the user can organise it as such using tags for example. The history contains the data initially uploaded (along with metadata such as the date and time of upload, the size, the way it was uploaded), all tools that have been used while the history is open (along with metadata such as date and time of launch, the run time, the energetic cost of the computation), all parameters, inputs and outputs (e.g. files, log, warnings) of each computations of tools. Users can decide to delete the results of given computations but the history keeps a record of these as well.

Title and abstract

Does the title clearly reflect the content of the article? Yes, No (please explain), I don't know

Does the abstract present the main findings of the study? Yes, No (please explain), I don't know

Introduction

Are the research questions/hypotheses/predictions clearly presented? Yes, No (please explain), I don't know

- As discussed above, it is difficult for the reader to initially figure out that Galaxy-Ecology is the focus of the paper.

The aim of the paper is now explained on lines 212-214, 239-243 and 559-569.

Galaxy-E is not the focus of the paper, with the shortened manuscript, especially sections about Galaxy, we hope it appears more clearly.

Does the introduction build on relevant research in the field? Yes, No (please explain), I don't know

Materials and methods

Are the methods and analyses sufficiently detailed to allow replication by other researchers? Yes, No (please explain), I don't know

Are the methods and statistical analyses appropriate and well described? Yes, No (please explain), I don't know

The paper is missing a detailed overview of the moving pieces of the tool.

As there is already a lot of existing material and publications on Galaxy, we cited and provided as much publications, tutorials and links in the manuscript to avoid further densification of the article. We only addressed technical details that served the message of the manuscript.

Results

In the case of negative results, is there a statistical power analysis (or an adequate Bayesian analysis or equivalence testing)? Yes, No (please explain), I don't know N/A

Are the results described and interpreted correctly? Yes, No (please explain), I don't know

Discussion

Have the authors appropriately emphasized the strengths and limitations of their study/theory/methods/argument? Yes, No (please explain), I don't know

There could be a little more comparison to other similar efforts to improve reproducibility and the limitations of what Galaxy-Ecology does.

We added a discussion and limitations section to further address these subjects.

Are the conclusions adequately supported by the results (without overstating the implications of the findings)? Yes, No (please explain), I don't know