



# How Reproducibility Will Accelerate Discovery Through Collaboration in Physio-Logging

Max F. Czapanskiy<sup>1\*</sup> and Roxanne S. Beltran<sup>2</sup>

<sup>1</sup>Hopkins Marine Station, Stanford University, Pacific Grove, CA, United States, <sup>2</sup>Department of Ecology and Evolutionary Biology, University of California Santa Cruz, Santa Cruz, CA, United States

## OPEN ACCESS

### Edited by:

Randall William Davis,  
Texas A&M University at Galveston,  
United States

### Reviewed by:

Gretchen Stahlman,  
Rutgers, The State University of New  
Jersey, United States  
Hazik Mohamed,  
Stellar Consulting Group, Singapore

### \*Correspondence:

Max F. Czapanskiy  
maxczapanskiy@gmail.com

### Specialty section:

This article was submitted to  
Physio-logging,  
a section of the journal  
Frontiers in Physiology

Received: 11 April 2022

Accepted: 16 June 2022

Published: 08 July 2022

### Citation:

Czapanskiy MF and Beltran RS (2022)  
How Reproducibility Will Accelerate  
Discovery Through Collaboration  
in Physio-Logging.  
Front. Physiol. 13:917976.  
doi: 10.3389/fphys.2022.917976

What new questions could ecophysio­logists answer if physio-logging research was fully reproducible? We argue that *technical debt* (computational hurdles resulting from prioritizing short-term goals over long-term sustainability) stemming from insufficient *cyberinfrastructure* (field-wide tools, standards, and norms for analyzing and sharing data) trapped physio-logging in a scientific silo. This debt stifles comparative biological analyses and impedes interdisciplinary research. Although physio-loggers (e.g., heart rate monitors and accelerometers) opened new avenues of research, the explosion of complex datasets exceeded ecophysiology's informatics capacity. Like many other scientific fields facing a deluge of complex data, ecophysio­logists now struggle to share their data and tools. Adapting to this new era requires a change in mindset, from "data as a noun" (e.g., traits, counts) to "data as a sentence", where measurements (nouns) are associated with transformations (verbs), parameters (adverbs), and metadata (adjectives). Computational reproducibility provides a framework for capturing the entire sentence. Though usually framed in terms of scientific integrity, reproducibility offers immediate benefits by promoting collaboration between individuals, groups, and entire fields. Rather than a tax on our productivity that benefits some nebulous greater good, reproducibility can accelerate the pace of discovery by removing obstacles and inviting a greater diversity of perspectives to advance science and society. In this article, we 1) describe the computational challenges facing physio-logging scientists and connect them to the concepts of *technical debt* and *cyberinfrastructure*, 2) demonstrate how other scientific fields overcame similar challenges by embracing computational reproducibility, and 3) present a framework to promote computational reproducibility in physio-logging, and bio-logging more generally.

**Keywords:** bio-logging, ecophysiology, ecoinformatics, cyberinfrastructure, technical debt

## INTRODUCTION

Ecophysiology, like many other scientific disciplines, is undergoing a technological revolution fueled by advances in both hardware and software. One technology dramatically advancing the research boundary is physio-logging, which enabled observations of animals' physiology in the field using animal-borne sensors (Fahlman et al., 2021; Hawkes et al., 2021). Ecophysio­logists suddenly find themselves with far more data, of increasing complexity, than ever before. Although these new data led to breakthroughs in answering individual questions, the lack of field-wide standards and tools largely siloed research groups, stifling collaboration and synthesis.

**TABLE 1** | Glossary of terms.

Term	Definition	References
Cyberinfrastructure	The collective interface between data collection and data analysis for a scientific field, including software, hardware, personnel, and shared practices	Atkins et al. (2003)
Technical debt	Short-term, sub-optimal choices in data and code that hamper future development without refactoring, such as missing documentation and bug-prone code	(Hinsen, 2015; Codabux et al., 2021; Vidoni, 2021)
Heterogeneous data	Combinations of data collected at different temporal scales and/or with different properties, for example multivariate time series (e.g., acceleration) with intermittent geospatial locations (e.g., GPS)	(Leinfelder et al., 2011; Michener and Jones, 2012)
Literate programming	A programming technique that combines code itself with descriptive text and outputs (figures, tables). R Markdown is an implementation of literate programming	(Knuth, 1984; Baumer and Udwin, 2015; Kery et al., 2018)
Data provenance	A record of the origin and processing steps that produced the data	Michener and Jones, (2012)

In this article, we introduce the concepts of *cyberinfrastructure* and *technical debt* (Table 1) as they apply to physio-logging research using a trio of researcher personas and the scientific challenges they face. Then, we describe how other scientific disciplines, such as astronomy, neuroscience, and molecular biology, have addressed these issues by embracing reproducibility as a guiding principle. Finally, we propose physio-logging cyberinfrastructure that promotes reproducibility and reduces barriers to collaboration.

## Barriers to Collaboration and Synthesis

In the following vignettes, three physio-logging researcher personas encounter common scientific challenges stemming from accumulated technical debt and insufficient cyberinfrastructure. We expect these challenges will be familiar to the reader and help connect their own experiences to the concepts in this article.

*Person A is a graduate student studying ecophysiology. They have their own physio-logging data and a second dataset contributed by a collaborating lab. Before conducting any meaningful analysis, Person A spends multiple weeks corresponding with another grad student in the collaborating lab to figure out how to merge the datasets due to mismatched variable names and units. Later in the analysis, Person A discovers that their own data recorded time in Coordinated universal Time and the collaborating lab used local time, requiring them to spend more time correcting the input data and re-running earlier steps. The analysis is not ready in time for their next committee meeting.*

*Person B is the Principal Investigator of a comparative physiology lab. The lab recently completed a field season and the trainees are working hard to process all the new physio-logging data. Meanwhile, Person B wants to know how their study system fits into the broader phylogenetic and morphological space, so they ask comparable physio-logging data from other labs. Despite a rich literature on this subject, only a small number of labs have the bandwidth and/or interest to contribute data for a synthesis study. As a result, Person B writes a paper with limited comparative power and they struggle to find funding for the next field season.*

*Person C is a data scientist at a government agency developing a new deep learning method for time series classification. They meet Person A at a café on campus and realize their physio-logging data*

*would be a perfect case study for the new method. Unfortunately, the dataset is not large enough on its own, so Person A promises to contact their collaborators to get more data. The data management issues encountered earlier by Person A present such substantial barriers to collating a sufficiently large dataset that both Persons lose interest and move on to other projects. Person C's method contributes to breakthrough advances in bio-medical research, but they never think about physio-logging again.*

The challenges encountered by the researcher personas illustrate how technical debt prevents collaboration and innovation. *Person A* is the early career researcher who works hands-on with bio-logging data, navigating a fragile computational ecosystem of inconsistently formatted data and poorly documented scripts. *Person B* is further along in their career and delegates computational tasks to achieve their scientific goals. *Person C* represents scientists in other fields who could make important contributions to physio-logging if they were given access and resources. The obstacles facing each persona are unique to their perspective and experience, but they all have the same consequence: narrowing the scope of scientific inquiry.

## Cyberinfrastructure and Technical Debt

*Cyberinfrastructure* is a layer in the technological stack driving modern science, situated between data production and discipline-specific analysis practices (Atkins et al., 2003). Translated into the physio-logging domain, “data production” includes electronic tags deployed on animals in the field. “Analysis practices” refers to both technical (e.g., statistical methods and packages) and social (e.g., norms about sharing data and code) aspects of physio-logging science. Cyberinfrastructure is the oft invisible middle layer that interfaces between data production and analysis practices. Efforts towards cyberinfrastructure in bio-logging generally include, for example, universal data standards (Campbell et al., 2016; Sequeira et al., 2021) and repositories with interfacing software (Kranstauber et al., 2011; Kays et al., 2022). The purpose of cyberinfrastructure is to “provide an effective and efficient platform for the empowerment of specific communities of researchers to innovate and eventually revolutionize what they do, how they do it, and who participates” (Atkins et al., 2003, p. 5).

In the absence of cyberinfrastructure, scientists are incentivized by funding, hiring, and promotion structures to choose rapid data collection and analysis over long-term technical sustainability. *Technical debt* is the lingering bug, the missing documentation, the “I swear it works on my machine” left behind in pursuit of scientific output (Hinsen, 2015). Up to a point, technical debt is not a problem; rather, it is the natural side effect of important scientific tasks like exploratory analyses and prototyping tools. But eventually technical debt creates obstacles for future work (Codabux et al., 2021; Vidoni, 2021). Removing those obstacles requires either 1) resources allocated specifically to auditing and fixing data and code (i.e., refactoring Adorf et al. (2019)) or 2) cyberinfrastructure that promotes best practices from the start. The first approach is unlikely to work at scale because existing incentive structures (funding, publications, hiring and promotion) do not prioritize refactoring existing data and code. Other scientific fields embraced the second approach, improved cyberinfrastructure, to address the same issues facing physio-logging. But before we investigate those efforts, we first describe in greater detail the technical debts in physio-logging.

## Where Did Our Technical Debt Come From

The underlying cause of the obstacles to collaboration and synthesis is a field-wide, collective technical debt incurred when an empirical/theoretical science (ecophysiology) embraced a data-intensive method (physio-logging) without simultaneously developing cyberinfrastructure. What was it about physio-logging that triggered ecophysiology’s technical debt?

Physio-logging changed the nature of ecophysiological data from simple, small tables to dense, *heterogeneous* (Table 1) mixtures of physiological, geospatial, and biomechanical time series (Harrison, 2021). Major developments in ecophysiology during the 20th and early 21st centuries were fueled by comparative analyses across morphology, phylogeny, geography, and other biological dimensions. Critically, these comparative studies were only possible because data could be combined from multiple individual studies. Consider the Metabolic Theory of Ecology, which emerged from scaling approaches that connected patterns and processes across the many levels of biological organization, from molecules to ecosystems (Brown et al., 2004). Three *Persons A* (Gillooly, Allen, and Savage—post-doctoral scholars at the time) and a *Person B* (Brown—an established biologist) assembled a synthetic dataset from dozens of earlier publications. In collaboration with a *Person C* (West—a theoretical physicist), the group articulated a theoretical framework for ecology from first principles. Though controversial, the Metabolic Theory of Ecology inspired a wave of innovative research (Lafferty et al., 2008; Burton et al., 2011; Gardner et al., 2011; Locey and Lennon, 2016) and amounted to a major advance in biological theory. Coincidentally, while Brown et al. were analyzing “simple” ecophysiological data, physio-logging data were growing in size and complexity. No longer simple recorders that measured individual data points (Kooyman, 1966; Goldbogen and Meir, 2014), by 2004 bio-loggers had evolved into high-

resolution, multisensory devices collecting gigabytes of heterogeneous data (Wilmers et al., 2015). A scientific community accustomed to simple, small datasets was suddenly presented with a profound informatics challenge—largely without the training, resources, and incentives to properly meet it. We have yet to develop adequate cyberinfrastructure, leaving a growing technical debt that impedes *Persons A-C*’s progress and limits discoveries using physio-logging tools.

## Reproducible Research Repays Technical Debt

As biologists, we are trained to think of data as observations. In other words, data are nouns: a cow’s mass, a salmon’s heart rate, a bird’s body temperature. However, physio-loggers record such vast quantities of complex, heterogeneous data that their interpretation *cannot be separated* from the computational methods used to process and analyze them. For example, consider a table of 50 Hz tri-axial accelerometer and magnetometer data. Even when visualized graphically, these data are largely meaningless to a human interpreter. But if we correct for the orientation of the tag relative to the animal’s body and account for the declination and inclination of the local magnetic field, then we can transform acceleration and magnetism into pitch, roll, and heading (Johnson and Tyack, 2003). Now the human interpreter can meaningfully visualize the animal’s fine-scale movements to identify physiologically relevant behaviors such as exercise (Williams et al., 2020), rest (Mitani et al., 2010), and escape (Williams et al., 2017). Physio-logging “data” are more than the final numbers; they are also the computational pipeline that transforms raw measurements into useful information. When it comes to physio-logging, data are not nouns; they are whole sentences composed of verbs (transformations), adjectives (metadata), and adverbs (parameters).

Computational reproducibility is the practice of writing the data’s entire sentence so that anyone can read it. Sharing our data this way accomplishes both “global” and “local” goals (Feinberg et al., 2020). The “global” goal is satisfying the scientific norm of reproducibility, providing transparency and integrity for our entire field. These ideals are the focus of the widely reported “reproducibility crisis” (Peng, 2015; Fanelli, 2018). But other data-intensive fields have embraced a “local” goal for computational reproducibility by reframing it in terms of collaboration and knowledge transfer. In this context, reproducibility is accomplished through computational best practices, such as re-usable code and proper documentation. This framing contains tangible solutions to the challenges experienced by *Persons A-C*. Sharing data as reproducible workflows (Cohen-Boulakia et al., 2017; Grüning et al., 2018; Wratten et al., 2021) solves the incompatibility issues that distracted *Person A* from doing good science and promotes collaboration within (*Person B*) and between (*Person C*) scientific disciplines. More importantly, it removes technical obstacles preventing broad, equitable access to our science. But a “data as a sentence”, reproducibility-focused mindset is too much to expect of individual ecophysiologicalists without adequate

cyberinfrastructure to support them. How have other fields provided tools, education, and other resources to their scientists?

## Cyberinfrastructure Examples From Other Fields

Most scientific fields are facing similar informatics challenges as ecophysiology; a few have developed cyberinfrastructure to promote sharing “data as a sentence”. What physio-logging is to ecophysiology, sky imaging is to astronomy, brain imaging is to neuroscience, and high-throughput sequencing is to molecular biology. The quantity and complexity of image and sequence data created technical debts in those fields as well, which were addressed through the coordinated development of cyberinfrastructure. Saliently, these fields used the “local” goal of computational reproducibility (collaboration and knowledge transfer) to motivate adoption of their respective cyberinfrastructures.

Astronomy’s cyberinfrastructure is represented by the aptly named Virtual Observatory (VO) (Quinn et al., 2004). There are currently 22 VOs around the world, such as the National Virtual Observatory in the United States and the Virtual Observatory of India. Each VO provides open access to sky imagery and other astronomical data in standardized formats agreed upon by the International Virtual Observatory Alliance (IVOA). In addition to data sharing, VOs provide open processing and analysis workflows, providing astronomers with invaluable tools for the “data as a sentence” mindset (Cui et al., 2020). The VO framework fuels innovative breakthroughs in astronomy, including the discovery of galaxies (Chilingarian et al., 2009).

In neuroscience, the proliferation of brain image data led to ad hoc, incompatible data curation practices, inspiring a standardization effort: the Brain Imaging Data Structure (BIDS). BIDS was developed to solve a problem familiar to physio-logging scientists: “Lack of consensus leads to misunderstanding and time wasted on rearranging data or rewriting scripts that expect particular file formats and organization, as well as a possible cause for errors.” (Gorgolewski et al., 2016, p. 2). In addition to a standardized data format, BIDS includes a software ecosystem for importing, validating, and processing imaging data (Gorgolewski et al., 2017). Like VOs, BIDS integrates data and code, facilitating “data as a sentence”.

The Human Genome Project published a draft of the human genome in 2001, representative of molecular biology’s pivot to big data (Lander et al., 2001). Three years later, Bioconductor emerged as a provider of data access, processing, and analysis tools (Gentleman et al., 2004; Huber et al., 2015). Like the VOs and BIDS, Bioconductor provides both data structures and computational methods. From the beginning, it was designed with reproducibility as an explicit goal, to promote collaboration in both sharing data and developing methods. As important new technologies emerge, such as high-throughput single-cell sequencing, the cyberinfrastructure provided by Bioconductor promotes collaborations between biologists, statisticians, and computer scientists to rapidly develop new methods for handling the influx of increasingly heterogeneous data (Amezquita et al., 2020). In turn, the cultural norm of

reproducibility incentivizes best practices, such as software documentation and validation, that facilitate adoption by researchers across the discipline.

VOs, BIDS, and Bioconductor exemplify the “data as a sentence” perspective by providing both data formats and computational tools for handling large quantities of heterogeneous data. Although there are data formats (Sequeira et al., 2021; Kays et al., 2022) and computational tools (Joo et al., 2020) for geospatial bio-logging data, the two sides have been developed independently without shared interfaces. Physio-logging lacks any such infrastructure.

## INTRODUCING `biologr`

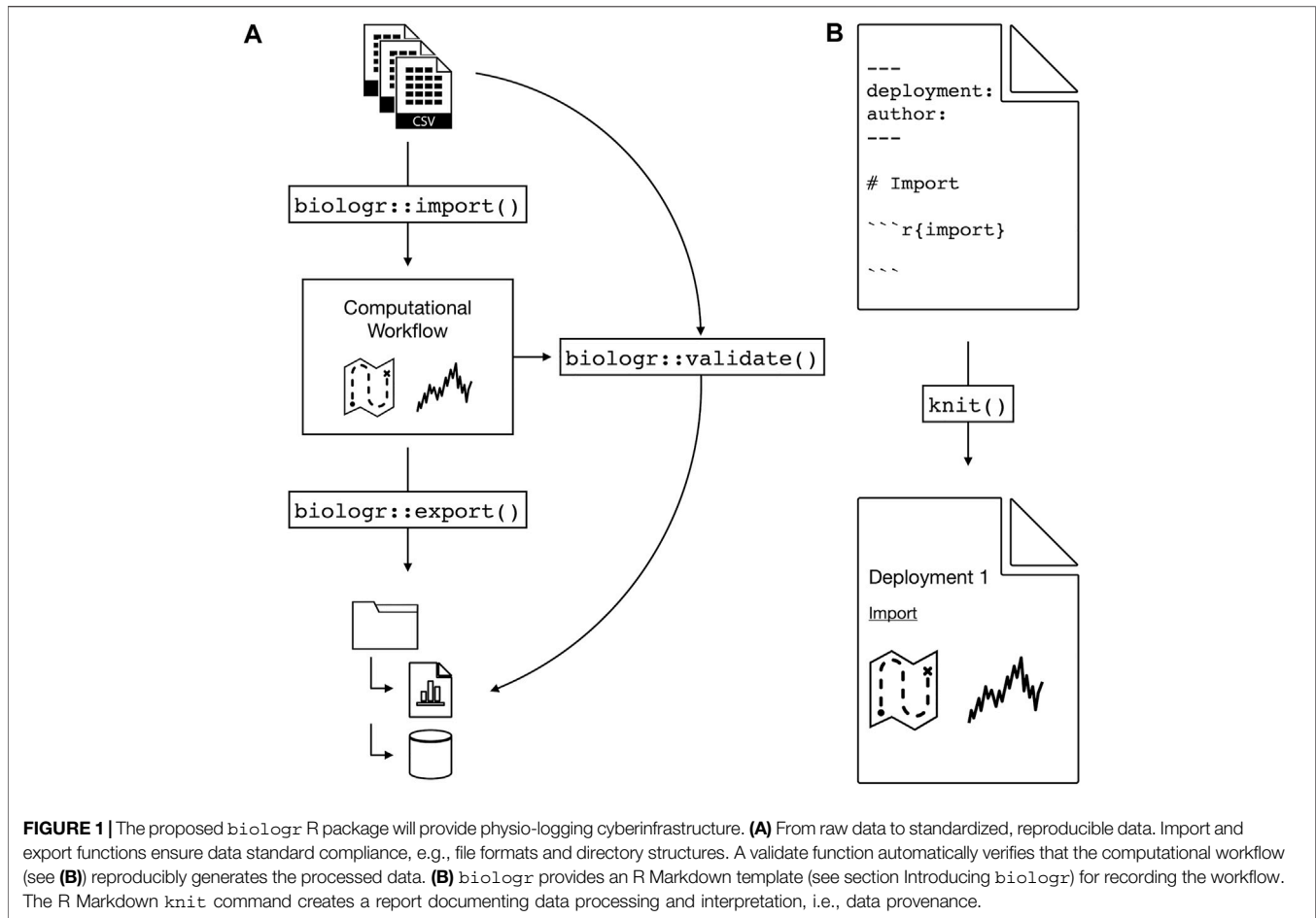
What would physio-logging cyberinfrastructure look like in practice? Following the lessons of the Virtual Observatory, BIDS, Bioconductor, and other successful efforts, it should facilitate reproducibility by integrating data formats directly with computational tools. Beginning with the data format, the nc-eTAG specification provides an extensible and efficient structure for storing bio-logging data (Tsontos et al., 2020). But recall the data alone are only the nouns; we must also capture the rest of the sentence. Existing workflows for processing bio-logging data are typically ad hoc and ephemeral, meaning they’re lost from the scientific record after publication. As a result, the data shared in repositories, whether they’re specialized for bio-logging like MoveBank or general purpose like Zenodo, are missing much of the story.

We propose an R package, `biologr`, that advances the goal of bio-logging cyberinfrastructure, supporting ecophysicologists and our colleagues in the broader bio-logging space. `biologr` interacts with the nc-eTAG data format and explicitly captures the entire workflow. By archiving both raw and processed data together with the workflow connecting them, `biologr` records the data’s entire sentence: nouns, verbs, and all (Figure 1).

What minimum set of functionalities would `biologr` require to make physio-logging data reproducible? In brief: data standards, import/export, validation, and workflow. For data standards, `biologr` relies on nc-eTAG for processed data and a standardized directory structure that organizes components for reproducibility. `biologr` needs both import and export functions to get the raw data into the nc-eTAG format and to populate the directory structure. Importantly, the export function saves both the data *and the computational workflow*. BIDS and other cyberinfrastructures have demonstrated that another critical function is automated validation (Gorgolewski et al., 2016). Computational reproducibility is built on standards, and without validation standardization quickly falls apart. So `biologr`’s automated validator checks both form (i.e., do the file formats and directory structures adhere to the standard?) and function (i.e., does the workflow reproducibly generate the processed data from the raw?). Import/export and validation both require a reproducible workflow, which is the component most responsible for reproducibility.

`biologr` captures the physio-logging computational workflow using *literate programming* (Table 1) (Knuth, 1984;





Kery et al., 2018). Literate programming is a technique that interweaves code, text, and output (e.g., figures and tables) - a combination that captures computational workflows in both human- and computer-readable formats. `bioloR` provides an R Markdown template, a literate programming implementation with broad adoption in the scientific community (Baumer and Udwin, 2015). Instead of an ad hoc, ephemeral script, the physio-logging researcher authors their workflow in R Markdown. Not only does this preserve computational details for automated validation and re-use/modification by other researchers, the output of the R Markdown file is itself a *provenance* report (Table 1), explicitly recording the transformations and interpretations in the data (Michener and Jones, 2012). This workflow approach promotes reproducibility by serving as the glue that binds together data standards, import/export, and validation.

## DISCUSSION

### A Vision for the Future of Physio-Logging

The goal of this article is to invite the physio-logging community to a conversation about the existing limitations, and potential advances, in our work relating to cyberinfrastructure and

reproducibility. The International Virtual Observatory Alliance, Brain Imaging Data Structure, and Bioconductor all have governing bodies with working groups to design, develop, and disseminate cyberinfrastructure for their fields. A working group within the International Bio-logging Society ([www.bio-logging.net](http://www.bio-logging.net)) could serve the same role for physio-logging, especially if ecoinformaticians were included in the process (Michener et al., 2012). The proposed `bioloR` package provides a starting point for those conversations. Below, we offer alternative vignettes to describe how our researcher personas would benefit from `bioloR`.

*Person A is a grad student studying ecophysiology. They have their own physio-logging data and a second dataset contributed by a collaborating lab. The two labs use slightly different methods to process their data, but both use `bioloR` for recording their workflows. Person A edits their collaborator's workflow to make the two datasets compatible and kicks off the reproducible workflow with a single click. They leave for lunch and come back to the lab to find two happily interoperable datasets. The analysis is done by the end of the week and their committee congratulates them on their progress.*

*Person B is the Principal Investigator of a comparative physiology lab. The lab recently completed a field season and the trainees are working hard to process all the new physio-logging data. Meanwhile, Person B wants to know how their study system fits into the broader*

phylogenetic and morphological picture. There is a rich literature on this subject, so dozens of archived datasets are available through a physio-logging data portal based on `biologr`. Person B downloads the data, runs a phylogenetically-informed scaling analysis, and writes a high-impact synthesis study. The publication opens exciting lines of collaborative inquiry and leads to a new multi-million dollar grant.

Person C is a post-doc in computer science developing a new deep learning method for time series classification. They meet Person A at a cafe on campus and realize their physio-logging data would be a perfect case study for the new method. Person C spends the afternoon reading `biologr` documentation and assembles a large physio-logging dataset from multiple studies. They demonstrate that their method accurately identifies physiological states from behavioral data and publish the method as a `biologr`-compatible R package, which is used by several physio-logging labs in their research.

## CONCLUSION

Physio-logging, and bio-logging more generally, gave biologists new tools for observing animals in their natural habitats with previously inconceivable detail. But concurrent with this great leap forward, we accumulated a technical debt that cost us the ability to easily share and synthesize our data. Ecophysiology became a data-intensive science without developing the cyberinfrastructure to handle large quantities of complex, heterogeneous data. `biologr` illustrates how we can develop tools to reproducibly process and archive physio-logging data. By embracing reproducibility, we can repay our technical debt and usher in a new era of collaboration while fostering best practices in a diverse next generation of physio-logging researchers. As individuals, labs, data repositories, and

## REFERENCES

- Adorf, C. S., Ramasubramani, V., Anderson, J. A., and Glotzer, S. C. (2019). How to Professionally Develop Reusable Scientific Software-And when Not to. *Comput. Sci. Eng.* 21, 66–79. doi:10.1109/MCSE.2018.2882355
- Amezquita, R. A., Lun, A. T. L., Becht, E., Carey, V. J., Carpp, L. N., Geistlinger, L., et al. (2020). Orchestrating Single-Cell Analysis with Bioconductor. *Nat. Methods* 17, 137–145. doi:10.1038/s41592-019-0654-x
- Atkins, D. E., Droegemeier, K., Feldman, S. I., Garcia-Molina, H., Klein, M., Messerschmitt, D., et al. (2003). *Revolutionizing Science and Engineering through Cyberinfrastructure, Report of the Blue-Ribbon Advisory Panel on Cyberinfrastructure*. Alexandria, VA: National Science Foundation.
- Baumer, B., and Udwin, D. (2015). R Markdown. *WIREs Comput. Stat.* 7, 167–177. doi:10.1002/wics.1348
- Brown, J. H., Gillooly, J. F., Allen, A. P., Savage, V. M., and West, G. B. (2004). Toward a Metabolic Theory of Ecology. *Ecology* 85, 1771–1789. doi:10.1890/03-9000
- Burton, T., Killen, S. S., Armstrong, J. D., and Metcalfe, N. B. (2011). What Causes Intraspecific Variation in Resting Metabolic Rate and what Are its Ecological Consequences? *Proc. R. Soc. B* 278, 3465–3473. doi:10.1098/rspb.2011.1778
- Campbell, H. A., Urbano, F., Davidson, S., Dettki, H., and Cagnacci, F. (2016). A Plea for Standards in Reporting Data Collected by Animal-Borne Electronic Devices. *Anim. Biotelemetry* 4, 1. doi:10.1186/s40317-015-0096-x

other stakeholders adopt shared standards, the possibilities for exploration and synthesis will grow exponentially.

## DATA AVAILABILITY STATEMENT

The original contributions presented in the study are included in the article/Supplementary Material, further inquiries can be directed to the corresponding author.

## AUTHOR CONTRIBUTIONS

Perspective conception: MC. Draft manuscript preparation: MC and RB.

## FUNDING

The work was funded in part by National Science Foundation award #2052497 to RB.

## ACKNOWLEDGMENTS

The authors thank Philippine Chambault, Millie Chapman, James Fahlbusch, William Oestreich, and Allison Payne for their constructive criticism of this manuscript. MC's perspective on the links between reproducibility and collaboration were greatly influenced by conversations with Susan Holmes' lab group and the Stanford Data Science Center for Open and REproducible Science seminar series. We are grateful for the insightful feedback from two reviewers.

- Chilingarian, I., Cayatte, V., Revaz, Y., Dodonov, S., Durand, D., Durret, F., et al. (2009). A Population of Compact Elliptical Galaxies Detected with the Virtual Observatory. *Science* 326, 1379–1382. doi:10.1126/science.1175930
- Codabux, Z., Vidoni, M., and Fard, F. H., 2021. Technical Debt in the Peer-Review Documentation of R Packages: a rOpenSci Case Study, Proceeding of the 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR). Presented at the 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR), June 2021, pp. 195–206. doi:10.1109/MSR52588.2021.00032
- Cohen-Boulakia, S., Belhajjame, K., Collin, O., Chopard, J., Froidevaux, C., Gaignard, A., et al. (2017). Scientific Workflows for Computational Reproducibility in the Life Sciences: Status, Challenges and Opportunities. *Future Gener. Comput. Syst.* 75, 284–298. doi:10.1016/j.future.2017.01.012
- Cui, C., Tao, Y., Li, C., Fan, D., Xiao, J., He, B., et al. (2020). Towards an Astronomical Science Platform: Experiences and Lessons Learned from Chinese Virtual Observatory. *Astronomy Comput.* 32, 100392. doi:10.1016/j.ascom.2020.100392
- Fahlman, A., Aoki, K., Bale, G., Brijs, J., Chon, K. H., Drummond, C. K., et al. (2021). The New Era of Physio-Logging and Their Grand Challenges. *Front. Physiol.* 12, 669158. doi:10.3389/fphys.2021.669158
- Fanelli, D. (2018). Is Science Really Facing a Reproducibility Crisis, and Do We Need it to? *Proc. Natl. Acad. Sci. U.S.A.* 115, 2628–2631. doi:10.1073/pnas.1708272114
- Feinberg, M., Sutherland, W., Nelson, S. B., Jarrahi, M. H., and Rajasekar, A. (2020). The New Reality of Reproducibility: The Role of Data Work in Scientific Research. *Proc. ACM Hum.-Comput. Interact.* 4, 0351–0352. doi:10.1145/3392840

- Gardner, J. L., Peters, A., Kearney, M. R., Joseph, L., and Heinsohn, R. (2011). Declining Body Size: a Third Universal Response to Warming? *Trends Ecol. Evol.* 26, 285–291. doi:10.1016/j.tree.2011.03.005
- Gentleman, R. C., Carey, V. J., Bates, D. M., Bolstad, B., Dettling, M., Dudoit, S., et al. (2004). Bioconductor: Open Software Development for Computational Biology and Bioinformatics. *Genome Biol.* 5, R80. doi:10.1186/gb-2004-5-10-r80
- Goldbogen, J. A., and Meir, J. U. (2014). The Device that Revolutionized Marine Organismal Biology. *J. Exp. Biol.* 217, 167–168. doi:10.1242/jeb.092189
- Gorgolewski, K. J., Auer, T., Calhoun, V. D., Craddock, R. C., Das, S., Duff, E. P., et al. (2016). The Brain Imaging Data Structure, a Format for Organizing and Describing Outputs of Neuroimaging Experiments. *Sci. Data* 3, 160044. doi:10.1038/sdata.2016.44
- Gorgolewski, K. J., Alfaro-Almagro, F., Auer, T., Bellec, P., Capotà, M., Chakravarty, M. M., et al. (2017). BIDS Apps: Improving Ease of Use, Accessibility, and Reproducibility of Neuroimaging Data Analysis Methods. *PLoS Comput. Biol.* 13, e1005209. doi:10.1371/journal.pcbi.1005209
- Grüning, B., Chilton, J., Köster, J., Dale, R., Soranzo, N., van den Beek, M., et al. (2018). Practical Computational Reproducibility in the Life Sciences. *Cell Syst.* 6, 631–635. doi:10.1016/j.cels.2018.03.014
- Harrison, X. A. (2021). A Brief Introduction to the Analysis of Time-Series Data from Biologging Studies. *Phil. Trans. R. Soc. B* 376, 20200227. doi:10.1098/rstb.2020.0227
- Hawkes, L. A., Fahlman, A., and Sato, K. (2021). Introduction to the Theme Issue: Measuring Physiology in Free-Living Animals. *Phil. Trans. R. Soc. B* 376, 20200210. doi:10.1098/rstb.2020.0210
- Hinsen, K. (2015). Technical Debt in Computational Science. *Comput. Sci. Eng.* 17, 103–107. doi:10.1109/MCSE.2015.113
- Huber, W., Carey, V. J., Gentleman, R., Anders, S., Carlson, M., Carvalho, B. S., et al. (2015). Orchestrating High-Throughput Genomic Analysis with Bioconductor. *Nat. Methods* 12, 115–121. doi:10.1038/nmeth.3252
- Johnson, M. P., and Tyack, P. L. (2003). A Digital Acoustic Recording Tag for Measuring the Response of Wild Marine Mammals to Sound. *IEEE J. Ocean. Eng.* 28, 3–12. doi:10.1109/OJEO.2002.808212
- Joo, R., Boone, M. E., Clay, T. A., Patrick, S. C., Clusella-Trullas, S., and Basille, M. (2020). Navigating through the R Packages for Movement. *J. Anim. Ecol.* 89, 248–267. doi:10.1111/1365-2656.13116
- Kays, R., Davidson, S. C., Berger, M., Bohrer, G., Fiedler, W., Flack, A., et al. (2022). The Movebank System for Studying Global Animal Movement and Demography. *Methods Ecol. Evol.* 13, 419–431. doi:10.1111/2041-210X.13767
- Kery, M. B., Radensky, M., Arya, M., John, B. E., and Myers, B. A. (2018). The Story in the Notebook. Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, CHI '18, Montréal, QC (New York, NY, USA: Association for Computing Machinery), 1–11. doi:10.1145/3173574.3173748
- Knuth, D. E. (1984). Literate Programming. *Comput. J.* 27, 97–111. doi:10.1093/comjnl/27.2.97
- Kooyman, G. L. (1966). Maximum Diving Capacities of the Weddell Seal, *Leptonychotes Weddelli*. *Science* 151, 1553–1554. doi:10.1126/science.151.3717.1553
- Kranstauber, B., Cameron, A., Weinzerl, R., Fountain, T., Tilak, S., Wikelski, M., et al. (2011). The Movebank Data Model for Animal Tracking. *Environ. Model. Softw.* 26, 834–835. doi:10.1016/j.envsoft.2010.12.005
- Lafferty, K. D., Allesina, S., Arim, M., Briggs, C. J., De Leo, G., Dobson, A. P., et al. (2008). Parasites in Food Webs: the Ultimate Missing Links. *Ecol. Lett.* 11, 533–546. doi:10.1111/j.1461-0248.2008.01174.x
- Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., et al. (2001). Initial Sequencing and Analysis of the Human Genome. *Nature* 409, 860–921. International Human Genome Sequencing Consortium, Whitehead Institute for Biomedical Research, C. for G.R., The Sanger Centre, Washington University Genome Sequencing Center, US DOE Joint Genome Institute, Baylor College of Medicine Human Genome Sequencing Center, RIKEN Genomic Sciences Center, Genoscope and CNRS UMR-8030; Department of Genome Analysis, I. of M.B., GTC Sequencing Center, Beijing Genomics Institute/Human Genome Center, Multimegabase Sequencing Center, T.I. for S.B., Stanford Genome Technology Center, University of Oklahoma's Advanced Center for Genome Technology, Max Planck Institute for Molecular Genetics, Cold Spring Harbor Laboratory, L.A.H.G.C., GBF—German Research Centre for Biotechnology, \*Genome Analysis Group (listed in alphabetical order, also includes individuals listed under other headings); Scientific management: National Human Genome Research Institute, U.N.I. of H., Stanford Human Genome Center, University of Washington Genome Center, Department of Molecular Biology, K.U.S. of M., University of Texas Southwestern Medical Center at Dallas, Office of Science, U.D. of E., The Wellcome Trust. doi:10.1038/35057062
- Leinfelder, B., Bowers, S., O'Brien, M., Jones, M. B., and Schildhauer, M. (2011). "Using Semantic Metadata for Discovery and Integration of Heterogeneous Ecological Data," in *Proceedings of the Environmental Information Management Conference (EIM 2011)*. Editors M. B. Jones and C. Gries, 92–97.
- Locey, K. J., and Lennon, J. T. (2016). Scaling Laws Predict Global Microbial Diversity. *Proc. Natl. Acad. Sci. U.S.A.* 113, 5970–5975. doi:10.1073/pnas.1521291113
- Michener, W. K., Allard, S., Budden, A., Cook, R. B., Douglass, K., Frame, M., et al. (2012). Participatory Design of DataONE-Enabling Cyberinfrastructure for the Biological and Environmental Sciences. *Ecol. Inf.* 11, 5–15. doi:10.1016/j.ecoinf.2011.08.007
- Michener, W. K., and Jones, M. B. (2012). Ecoinformatics: Supporting Ecology as a Data-Intensive Science. *Trends Ecol. Evol.* 27, 85–93. doi:10.1016/j.tree.2011.11.016
- Mitani, Y., Andrews, R. D., Sato, K., Kato, A., Naito, Y., and Costa, D. P. (2010). Three-dimensional Resting Behaviour of Northern Elephant Seals: Drifting like a Falling Leaf. *Biol. Lett.* 6, 163–166. doi:10.1098/rsbl.2009.0719
- Peng, R. (2015). The Reproducibility Crisis in Science: A Statistical Counterattack. *Significance* 12, 30–32. doi:10.1111/j.1740-9713.2015.00827.x
- Quinn, P. J., Barnes, D. G., Csabai, I., Cui, C., Genova, F., Hanisch, B., et al. (2004). "The International Virtual Observatory Alliance: Recent Technical Developments and the Road Ahead," in *Proceedings of SPIE*. Vol. 5493, 137. doi:10.1117/12.551247
- Sequeira, A. M. M., O'Toole, M., Keates, T. R., McDonnell, L. H., Braun, C. D., Hoenner, X., et al. (2021). A Standardisation Framework for Bio-logging Data to Advance Ecological Research and Conservation. *Methods Ecol. Evol.* 12, 996–1007. doi:10.1111/2041-210X.13593
- Tsontos, V. M., Chihin, L., and Arms, S. (2020). *NASA-OIIP netCDF Templates for Electronic Tagging Data: The Ne-TAG File Format and Metadata Specification (Version 1.0)*. [WWW Document]. figshare. doi:10.6084/m9.figshare.10159820.v1
- Vidoni, M. (2021). Evaluating Unit Testing Practices in R Packages. Proceeding of the 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE). Presented at the 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), May 2021, Madrid, ES. IEEE, 1523–1534. doi:10.1109/ICSE43902.2021.00136
- Williams, T. M., Blackwell, S. B., Richter, B., Sinding, M.-H. S., and Heide-Jørgensen, M. P. (2017). Paradoxical Escape Responses by Narwhals ( *Monodon Monoceros* ). *Science* 358, 1328–1331. doi:10.1126/science.aao2740
- Williams, H. J., Shepard, E. L. C., Holton, M. D., Alarcón, P. A. E., Wilson, R. P., and Lambertucci, S. A. (2020). Physical Limits of Flight Performance in the Heaviest Soaring Bird. *Proc. Natl. Acad. Sci. U.S.A.* 117, 17884–17890. doi:10.1073/pnas.1907360117
- Wilmers, C. C., Nickel, B., Bryce, C. M., Smith, J. A., Wheat, R. E., and Yovovich, V. (2015). The Golden Age of Bio-Logging: How Animal-Borne Sensors Are Advancing the Frontiers of Ecology. *Ecology* 96, 1741–1753. doi:10.1890/14-1401.1
- Wratten, L., Wilm, A., and Göke, J. (2021). Reproducible, Scalable, and Shareable Analysis Pipelines with Bioinformatics Workflow Managers. *Nat. Methods* 18, 1161–1168. doi:10.1038/s41592-021-01254-9

**Conflict of Interest:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

**Publisher's Note:** All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Copyright © 2022 Czapanskiy and Beltran. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.