

***NATURE ECOLOGY AND EVOLUTION EMBARGO***  
***– MAY 23, 2017, 10:00 EDT –***

**Our path to better science in less time using open data  
science tools**

Julia S. Stewart Lowndes<sup>1\*</sup>, Benjamin D. Best<sup>2</sup>, Courtney Scarborough<sup>1</sup>, Jamie C.  
Afflerbach<sup>1</sup>, Melanie R. Frazier<sup>1</sup>, Casey C. O'Hara<sup>1</sup>, Ning Jiang<sup>1</sup>, Benjamin S.  
Halpern<sup>1,3,4</sup>

<sup>1</sup>National Center for Ecological Analysis and Synthesis, University of California at  
Santa Barbara, Santa Barbara, CA, United States

<sup>2</sup>EcoQuants.com, Santa Barbara, CA, United States

<sup>3</sup>Bren School for Environmental Science and Management, University of California,  
Santa Barbara, CA, United States

<sup>4</sup>Silwood Park Campus, Imperial College London, Ascot, United Kingdom

\*corresponding author: lowndes@nceas.ucsb.edu

## 15 Preface

16 Reproducibility has long been a tenet of science but has been challenging to achieve  
17 — we learned this the hard way when our old approaches proved inadequate to  
18 efficiently reproduce our own work. Here we describe how several free software  
19 tools have fundamentally upgraded our approach to collaborative research, making  
20 our entire workflow more transparent and streamlined. By describing specific tools  
21 and how we incrementally began using them for the Ocean Health Index project, we  
22 hope to encourage others in the scientific community to do the same — so we can all  
23 produce better science in less time.

24

## 25 Keywords

26 collaboration, data science, Ocean Health Index, open science, reproducibility,  
27 transparency

28

## 29 **Scientists need data science**

30 Science, now more than ever, demands reproducibility, collaboration, and effective  
31 communication to strengthen public trust and effectively inform policy. Recent high-  
32 profile difficulties in reproducing and repeating scientific studies have put the  
33 spotlight on psychology and cancer biology<sup>1-3</sup>, but it is widely acknowledged that  
34 reproducibility challenges persist across scientific disciplines<sup>4-6</sup>. Environmental  
35 scientists face potentially unique challenges in achieving goals of transparency and  
36 reproducibility because they rely on vast amounts of data spanning natural,  
37 economic, and social sciences that create semantic and synthesis issues exceeding  
38 those for most other disciplines<sup>7-9</sup>. Furthermore, proposed environmental solutions  
39 can be complex, controversial, and resource intensive, increasing the need for  
40 scientists to work transparently and efficiently with data to foster understanding  
41 and trust.

42 Environmental scientists are expected to work effectively with ever-increasing  
43 quantities of highly heterogeneous data even though they are seldom formally  
44 trained to do so<sup>10-14</sup>. This was recently highlighted by a survey of 704 US National  
45 Science Foundation principle investigators in the biological sciences that found  
46 training in data skills to be the largest unmet need<sup>15</sup>. Without training, scientists  
47 tend to develop their own bespoke workarounds to keep pace, but with this comes  
48 wasted time struggling to create their own conventions for managing, wrangling,  
49 and versioning data. If done haphazardly or without a clear protocol, these efforts  
50 are likely to result in work that is not reproducible — by the scientist's own 'future

self' or by anyone else<sup>12</sup>. As a team of environmental scientists tasked with reproducing our own science annually, we experienced this struggle first-hand. When we began our project, we worked with data in the same way as we always had, taking extra care to make our methods reproducible for planned future re-use. But when we began to reproduce our workflow a second time and repeat our methods with updated data, we found our approaches to reproducibility were insufficient. However, by borrowing philosophies, tools, and workflows primarily created for software development, we have been able to dramatically improve the ability for ourselves and others to reproduce our science, while also reducing the time involved to do so: the result is better science in less time (Fig. 1).

Here we share a tangible narrative of our transformation to better science in less time — meaning more transparent, reproducible, collaborative, and openly shared and communicated science — with an aim of inspiring others. Our story is only one potential path because there are many ways to upgrade scientific practices — whether collaborating only with your 'future self' or as a team — and they depend on the shared commitment of individuals, institutions, and publishers<sup>6,16,17</sup>. We do not review the important, ongoing work regarding data management architecture and archiving<sup>8,18</sup>, workflows<sup>11,19–21</sup>, sharing and publishing data<sup>22–25</sup> and code<sup>25–27</sup>, or how to tackle reproducibility and openness in science<sup>28–32</sup>. Instead, we focus on our experience, because it required changing the way we had always worked, which was extraordinarily intimidating. We give concrete examples of how we use tools and practices from data science, the discipline of 'turning raw data into understanding'<sup>33</sup>. It was out of necessity that we began to engage in data science,

which we did incrementally by introducing new tools, learning new skills, and creating deliberate workflows — all while maintaining annual deadlines. Through our work with academics, governments, and non-profit groups around the world, we have seen that the need to improve practices is common if not ubiquitous. In this narrative we describe specific software tools, why we use them, how we use them in our workflow, and how we work openly as a collaborative team. In doing so we underscore two key lessons we learned that we hope encourage others to incorporate these practices into their own research. The first is that powerful tools exist and are freely available to use; the barriers to entry seem to be exposure to relevant tools and building confidence using them. The second is that engagement may best be approached as an evolution rather than as a revolution that may never come.

## **Improving reproducibility and collaboration**

### **From then to now**

The Ocean Health Index (OHI) operates at the interface of data-intensive marine science, coastal management and policy, and now, data science<sup>34,35</sup>. It is a scientific framework to quantify ocean-derived benefits to humans and to help inform sustainable ocean management using the best available information<sup>36,37</sup>. Assessments using the OHI framework require synthesising heterogeneous data from nearly one hundred different sources, ranging from categorical tabular data to high-resolution rasters. Methods must be reproducible, so that others can produce

the same results, and also repeatable, so that newly available data can be incorporated in subsequent assessments. Repeated assessments using the same methods enable quantifiable comparison of changes in ocean health through time, which can be used to inform policy and track progress<sup>34</sup>.

Using the OHI framework, we lead annual global assessments of 220 coastal nations and territories, completing our first assessment in 2012<sup>36</sup>. Despite our best efforts, we struggled to efficiently repeat our own work during the second assessment in 2013 because of our approaches to data preparation<sup>37</sup>. Data preparation is a critical aspect of making science reproducible but is seldom explicitly reported in research publications; we thought we had documented our methods sufficiently in 130-pages of published supplemental materials<sup>36</sup>, but we had not.

However, by adopting data science principles and freely available tools that we describe below, we began building an OHI ‘Toolbox’ and fundamentally changed our approach to science (Figure 1). The OHI Toolbox provides a file structure, data, code, and instruction, operates across computer operating systems, and is shared online for free so that anyone can begin building directly from previous OHI assessments without reinventing the wheel<sup>34</sup>. While these changes required an investment of our team’s time to learn and develop the necessary skills, the pay-off has been substantial. Most significantly we are now able to share and extend our workflow with a growing community of government, non-profit, and academic collaborations around the world that use the OHI for science-driven marine management. There are currently two dozen OHI assessments underway, most of which are led by

independent groups<sup>34</sup>, and the Toolbox has helped lower the barriers to entry. Further, our own team has just released the fifth annual global OHI assessment<sup>38</sup> and continues to lead assessments at smaller spatial scales, including the Northeastern United States, where the OHI is included in President Obama's first Ocean Plan<sup>39</sup>.

## **We thought we were doing reproducible science**

For the first global OHI assessment in 2012 we employed an approach to reproducibility that is standard to our field, which focused on scientific methods, not data science methods<sup>36</sup>. Data from nearly one hundred sources were prepared manually — i.e. without coding, typically in Microsoft Excel — which included organising, transforming, rescaling, gap-filling, and formatting data. Processing decisions were documented primarily within the Excel files themselves, emails, and Microsoft Word documents. We programmatically coded models and meticulously documented their development, (resulting in the 130-page supplemental materials)<sup>36</sup>, and upon publication, we also made the model inputs (i.e., prepared data and metadata) freely available to download. This level of documentation and transparency is beyond the norm for environmental science<sup>16,40</sup>.

We also worked collaboratively in the same ways we always had. Our team included scientists and analysts with diverse skill sets and disciplines, and we had distinct, domain-specific roles assigned to scientists and to a single analytical programmer. Scientists were responsible for developing the models conceptually, preparing data, and interpreting modeled results, and the programmer was responsible for coding

the models. We communicated and shared files frequently, with long, often-forwarded, and vaguely-titled email chains (e.g. Re: Fwd: data question) with manually versioned data files (e.g. data\_final\_updated2.xls). All team members were responsible for organising those files with their own conventions on their local computers. Final versions of prepared files were stored on the servers and used in models, but records of the data processing itself were scattered.

Upon beginning the second annual assessment in 2013, we realised that our approach was insufficient since it took too much time and relied heavily on individuals' data organisation, email chains, and memory — particularly problematic as original team members moved on and new team members joined. We quickly realised we needed a nimble and robust approach to sharing data, methods, and results within and outside our team — we needed to completely upgrade our workflow.

## Actually doing reproducible science

As we began the second global OHI assessment in 2013 we faced challenges across three main fronts: 1) *reproducibility*, including transparency and repeatability, particularly in data preparation; 2) *collaboration*, including team record keeping and internal collaboration; and 3) *communication* with scientific and broader communities. We knew that environmental scientists are increasingly using R<sup>11</sup> because it is free, cross-platform, and open source, and also because of the training and support provided by developers<sup>33</sup> and independent groups<sup>12,41</sup> alike. We decided to base our work in R<sup>42</sup> and RStudio<sup>43</sup> for coding and visualisation, Git<sup>44</sup> for



version control, GitHub<sup>45</sup> for collaboration, and a combination of GitHub and RStudio for organisation, documentation, project management, online publishing, distribution, and communication (Table 1). These tools can help scientists organise, document, version, and easily share data and methods, thus not only increasing reproducibility but also reducing the amount of time involved to do so<sup>14,46,47</sup>. Many available tools are free so long as work is shared publicly online, which enables open science, defined by Hampton et al.<sup>40</sup> as "the concept of transparency at all stages of the research process, coupled with free and open access to data, code, and papers". When integrated into the scientific process, data science tools that enable open science — let's call them "open data science" tools — can help realise reproducibility in collaborative scientific research<sup>6,16,40,48,49</sup>.

Open data science tools helped us upgrade our approach to reproducible, collaborative, and transparent science, but they did require a substantial investment to learn, which we did incrementally over time (Figure 1; Box 1). Previous to this evolution, most team members with any coding experience — not necessarily in R — had learned just enough to accomplish whatever task had been before them using their own unique conventions. Given the complexity of the OHI project, we needed to learn to code collaboratively and incorporate best<sup>50,51</sup> or good enough practices<sup>12,52</sup> into our coding, so that our methods could be co-developed and vetted by multiple team members. Using a version control system not only improved our file and data management, but allowed individuals to feel less inhibited about their coding contributions, since files could always be reverted back to previous versions if there were problems. We built confidence using these tools by sharing our

imperfect code, discussing our challenges, and learning as a team. These tools quickly became the keystone of how we work, and have overhauled our approach to science, perhaps as much as email did in decades prior. They have changed the way we think about science and about what is possible. The following describes how we have been using open data science practices and tools to overcome the biggest challenges we encountered to reproducibility, collaboration, and communication.

## Reproducibility

### Data preparation - coding and documenting

Our first priority was to code all data preparation, create a standard format for final data layers, and do so using a single programmatic language, R<sup>42</sup>. Code enables us to reproduce the full process of data preparation, from data download to final model inputs<sup>37,53</sup>, and a single language makes it more practical for our team to learn and contribute collaboratively. We code in R and use RStudio<sup>43</sup> to power our workflow because it has a user-friendly interface and built-in tools useful for coders of all skill levels, and, importantly, it can be configured with Git to directly sync with GitHub online (See Collaboration section). We have succeeded in transitioning to R as our primary coding language for data preparation, including for spatial data, although some operations still require additional languages and tools such as ArcGIS, QGIS, and Python<sup>54–56</sup>.

All our code is underpinned by the principles of tidy data, the grammar of data manipulation, and the tidyverse R packages developed by Wickham<sup>33,57–59</sup>. This

205 deliberate philosophy for thinking about data helped bridge our scientific questions  
206 with the data processing required to get there, and the readability and conciseness  
207 of `tidyverse` operations makes our data analysis read more as a story arc.  
208 Operations require less syntax — which can mean fewer potential errors that are  
209 easier to identify — and they can be chained together, minimising intermediate  
210 steps and data objects that can cause clutter and confusion<sup>33,60</sup>. `tidyverse` tools for  
211 wrangling data have expedited our transformation as coders and made R less  
212 intimidating to learn. We heavily rely on a few packages for data wrangling and  
213 visualisation that are bundled in the `tidyverse` package<sup>58,59</sup> — particularly `dplyr`,  
214 `tidyr`, and `ggplot2` — as well as accompanying books, cheatsheets, and archived  
215 webinars (Box 1).

216 We keep detailed documentation describing metadata (e.g., source, date of access,  
217 links) and data processing decisions — trying to capture not only the processing we  
218 decided to do, but what we decided against. We started with small plain text files  
219 accompanying each R file, but have transitioned to documenting with R  
220 Markdown<sup>61,62</sup> because it combines plain text and executable chunks of R code  
221 within the same file and serves as a living lab notebook. Every time R Markdown  
222 output files are regenerated the R code is rerun so the text and figures will also be  
223 regenerated and reflect any updates to the code or underlying data. R Markdown  
224 files increase our reproducibility and efficiency by streamlining documentation and  
225 eliminating the need to constantly paste updated figures into reports as they are  
226 developed.

## Modeling – R functions and packages

Once the data are prepared, we develop assessment-specific models to calculate OHI scores. Models were originally coded in multiple languages to accommodate disparate data types and formatting. By standardising our approach to data preparation and final data layer format, we have been able to translate all models into R. In addition to assessment-specific models, the OHI framework includes core analytical operations that are used by all OHI assessments<sup>34</sup>, and thus we created an R package called `ohicore`<sup>63</sup>, which was greatly facilitated by the `devtools` and `roxygen2` packages<sup>64–66</sup>. The `ohicore` package is maintained in and installed from a dedicated GitHub repository — using `devtools::install_github('ohi-science/ohicore')` — from any computer with R and an internet connection, enabling groups leading independent OHI assessments to use it for their own work<sup>34</sup>.

## Version control

We use Git<sup>44</sup> as a version control system. Version control systems track changes within files and allow you to examine or rewind to previous versions. This saves time that would otherwise be spent duplicating, renaming, and organising files to preserve past versions. It also makes folders easier to navigate since they are no longer overcrowded with multiple files suffixed with dates or initials (e.g., `final_JL-2012-02-26.csv`)<sup>67–69</sup>. Once Git is configured on each team member's machine, they work as before but frequently commit to saving a snapshot of their

files, along with a human-readable "commit message"<sup>67,68</sup>. Any line modified in a file tracked by Git will then be attributed to that user.

We interface with Git primarily through RStudio, using the command line for infrequently encountered tasks. Using RStudio to interact with Git was key for our team's uptake of a version control system, since the command line can be an intimidating hurdle or even a barrier for beginners to get onboard with using version control. We were less resistant because we could use a familiar interface, and as we gained fluency in Git's operations through RStudio we translated that confidence to the command line.

## Organisation

Our team developed conventions to standardise the structure and names of files to improve consistency and organisation. Along with the GitHub workflow (see Collaboration section below), having a structured approach to file organisation and naming has helped those within and outside our team navigate our methods more easily. We organise parts of the project in folders that are both RStudio "projects" and GitHub "repositories", which has also helped us collaborate using shared conventions rather than each team member spending time duplicating and organising files.

## 266    **Collaboration within our team**

### 267    **Coding collaboratively**

268    We transitioned from a team of distinct roles (scientists-and-programmer) to  
269    becoming a team with overlapping skill sets (scientists-as-programmers, or simply,  
270    data scientists). Having both environmental expertise and coding skills in the same  
271    person increases project efficiency, enables us to vet code as a team, and reduces the  
272    bottleneck of relying on a single programmer. We, like Duhigg<sup>70</sup>, have found that  
273    "groups tend to innovate faster, see mistakes more quickly and find better solutions  
274    to problems". Developing these skills and creating the team culture around them  
275    requires leadership with the understanding that fostering more efficient and  
276    productive scientists is worth the long-term investment. Our team had the freedom  
277    to experiment with available tools and their value was recognised with a  
278    commitment that we, as a team, would adopt and pursue these methods further. In  
279    addition to supportive leadership, having a "champion" with experience of how  
280    tools can be introduced over time and interoperate can expedite the process, but is  
281    not the only path (Box 2). Taking the time to experiment and invest in learning data  
282    science principles, tools, and skills enabled our team to establish a system of best  
283    practices for developing, using, and teaching the OHI Toolbox.

### 284    **Our (simplified) GitHub workflow**

285    GitHub is one of many web-based platforms that enables files tracked with Git to be  
286    collaboratively shared online so contributors can keep their work

287 synchronised<sup>45,68,69</sup>, and it is increasingly being adopted by scientific communities  
288 for project management<sup>71</sup>. Versioned files are synced online with GitHub similar to  
289 the way Dropbox operates, except syncs require a committed, human-readable  
290 message and reflect deliberate snapshots of changes made that are attributed to the  
291 user, line-by-line, through time. Built for large, distributed teams of software  
292 developers, GitHub provides many features that we as a scientific team, new to data  
293 science, do not immediately need, and thus we mostly ignore features such as  
294 branching, forking, and pull requests. Our team uses a simplified GitHub workflow  
295 whereby all members have administrative privileges to the repositories within our  
296 ohi-science organisation. Each team member is able to sync their local work to  
297 GitHub.com, making it easier to attribute contribution, as well as identify to whom  
298 to direct questions.

299 GitHub is now central to many facets of our collaboration as a team and with other  
300 communities — we use it along with screensharing to teach and troubleshoot with  
301 groups leading independent OHI assessments, as well as to communicate our  
302 ongoing work and final results (see Communication section). Now there are very  
303 few files emailed back and forth within our team since we all have access to all  
304 repositories within the ohi-science organisation, and can navigate to and edit  
305 whatever we need. Additionally, these organised files are always found with the  
306 same file path, whether on GitHub.com or on someone's local computer; this, along  
307 with RStudio .Rproj files, eases the file path problems that can plague collaborative  
308 coding and frustrate new coders.

## 309 **Internal communication**

310 We use a feature of GitHub called 'Issues' in place of email for discussions about data  
311 preparation and analysis. We use Issues in a separate private repository to keep our  
312 conversations private but our work public. All team members can see and  
313 contribute to all conversations, which are a record of all our decisions and  
314 discussions across the project and are searchable in a single place. Team members  
315 can communicate clearly by linking to specific lines of code in current or past  
316 versions of specific files since they are stored on GitHub and thus have a URL, as  
317 well as paste images and screenshots, link to other websites, and send an email to  
318 specific team members directly by mentioning their GitHub username. In addition to  
319 discussing analytical options, we use Issues to track ongoing tasks, tricks we have  
320 learned, and future ideas. Issues provide a written reference of institutional memory  
321 so new team members can get up to speed more easily. Most importantly, GitHub  
322 Issues have helped us move past the never-ending forwarded email chains and  
323 instead to conversations available to any current or future team member.

## 324 **Communication outside the project**

### 325 **Sharing data and code**

326 Our code is online in GitHub repositories, publicly available for any researcher or  
327 interested person to see and access ([github.com/ohi-science](https://github.com/ohi-science)). As we work, GitHub  
328 renders code, text, images, and tabular and spatial data and displays differences  
329 between versions, essentially creating webpages that can be easily shared with



collaborators, whether or not they use GitHub. Additionally, we create ‘Releases’ for each global assessment<sup>36,37</sup> so the code and data we use for peer-reviewed publication are preserved while we continue our work (<https://github.com/OHI-Science/ohi-global/releases>).

## Sharing methods and instruction

We use R Markdown not only for data preparation but also for broader communication. R Markdown files can be generated into a wide variety of formatted outputs, including PDFs, slides, Microsoft Word documents, HTML files, books, or full websites<sup>61,62</sup>. These can all be published online for free through GitHub using the same RStudio-GitHub workflow that we use for our analyses, which has made communication an ongoing part of our work, instead of a final step in completed analyses.

We built a website using GitHub and RStudio publishing tools: [ohi-science.org](https://ohi-science.org). Team members can update content directly, and using the same workflow makes it easier for us to keep it current. Our website is intended for scientists interested in our methods as well as those leading their own assessments<sup>34</sup>. Thus, the website provides scientific methods, publications, data, and code, as well as instruction, news, blog posts, and a map displaying where all ongoing OHI assessments are taking place so that groups can learn directly from and build off of each other's code. [ohi-science.org](https://ohi-science.org) provides technical information to complement [oceanhealthindex.org](https://oceanhealthindex.org), our overview website intended for more general audiences.

## Meeting scientists where they are

We are environmental scientists whose impetus for upgrading approaches to collaborative, data-intensive science was driven by our great difficulty reproducing our own methods. Many researchers do not attempt to reproduce their own work<sup>17,72</sup> — ourselves included before 2013 — and thus may not realise that there could be reproducibility issues in their own approaches. But they can likely identify inefficiencies. Integrating open data science practices and tools into science can save time, while also improving reproducibility for our most important collaborator: our ‘future selves’. We have found this as individuals and as a team: We could not be as productive<sup>34,35</sup> without open data science practices and tools. We would also not be able to efficiently share and communicate our work while it is ongoing rather than only post-publication, which is particularly important for bridging science and policy. As environmental scientists who are still learning, we hope sharing our experiences will empower other scientists to upgrade their own approaches, helping further shift the scientific culture to value transparency and openness as a benefit to all instead of as a vulnerability<sup>16,40,48</sup>.

From our own experience and from teaching other academic, non-profit, and government groups through the Ocean Health Index project<sup>34</sup>, we find that the main barriers to engagement boil down to exposure and confidence: first knowing which tools exist that can be directly useful to one's research, and then having the confidence to develop the skills to use them. These two points are simple but critical. We are among the many environmental scientists who were never formally

373 trained to work deliberately with data. Thus, we were unaware of how significantly  
374 open data science tools could directly benefit our research<sup>11,73</sup>, and upon learning  
375 about them we were hesitant, or even resistant, to engage. However, we were able  
376 to develop confidence in large part because of the open, inclusive, and encouraging  
377 online developer community that builds tools and creates tutorials that meet  
378 scientists where they are (Box 1, Box 2). It takes motivation, patience, diligence, and  
379 time to overcome the conceptual and technical challenges involved in developing  
380 computing skills but resources are available to help scientists get started<sup>11,51,73</sup>.  
381 Coding is "as important to modern scientific research as telescopes and test  
382 tubes"<sup>50</sup>, but it is critical to "dispel the misconception that these skills are intuitive,  
383 obvious, or in any way inherent"<sup>41</sup>.

384 There is ongoing and important work by the informatics community on the  
385 architecture and systems for data management and archiving<sup>7,8,18,74</sup>, as well as  
386 efforts to enable scientists to publish the code that they do have<sup>26,31,52</sup>. This work is  
387 critical, but comes with the *a priori* assumption that scientists are already thinking  
388 about data and coding in a way that they would seek out further resources. In  
389 reality, this is not always the case, and without visible examples of how to use these  
390 tools within their scientific fields, common stumbling blocks will be continually  
391 combatted with individual workarounds instead of addressed with intention. These  
392 workarounds can greatly delay focusing on actual scientific research, particularly  
393 when scientific questions that may not yet have answers — e.g., how the behavior of  
394 X changes with Y — are conflated with data science questions that have many  
395 existing answers — e.g., how to operate on only criteria X and Y.

Scientific advancement comes from building off the past work of others; scientists can also embrace this principle for using software tools to tackle some of the challenges encountered in modern scientific research. In a recent survey in *Nature*, 90% of the 1,500 respondents across scientific fields agreed that there was a reproducibility crisis in science, and one third of the respondents reported not having their own "established procedures for reproducibility"<sup>4</sup>. While reproducibility means distinct things within the protocols of each sub-discipline or specialty, underpinning reproducibility across all disciplines in modern science is working effectively and collaboratively with data, including wrangling, formatting, and other tasks that can take 50-80% of a data scientist's time<sup>75</sup>. While reaching full reproducibility is extremely difficult<sup>5,76</sup>, incrementally incorporating open data science practices and tools into scientific workflows has the potential to alleviate many of the troubles plaguing science, including collaboration and preserving institutional memory<sup>12</sup>. Further, sharing openly is fundamental to truly expediting scientific progress because others can build directly off previous work if well-documented, re-usable code are available<sup>16,47,48,77</sup>. Until quite recently, making research open required a great deal of extra work for researchers and was less likely to be done. Now, with available tools, the benefits of openness can be a byproduct of time-saving efficiencies, because tools that reduce data headaches also result in science that is more transparent, reproducible, collaborative, and freely accessible to others.

Ecologists and environmental scientists arguably have a heightened responsibility for transparency and openness, as data products provide important snapshots of

419 systems that may be forever altered due to climate change and other human  
420 pressures<sup>16,18</sup>. There is particular urgency for efficiency and transparency, as well as  
421 opportunity to democratise science in fields that operate at the interface of science  
422 and policy. Individuals play an important part by promoting good practices and  
423 creating supportive communities<sup>16,41,48</sup>. But it is also critical for the broader science  
424 community to build a culture where openness and reproducibility are valued,  
425 formally taught, and practiced, where we all agree that they are worth the  
426 investment.

427

## References

1. Baker, M. Over half of psychology studies fail reproducibility test. *Nature* (2015). doi:[10.1038/nature.2015.18248](https://doi.org/10.1038/nature.2015.18248)
2. Baker, M. & Dolgin, E. Cancer reproducibility project releases first results. *Nature News* **541**, 269 (2017). doi:10.1038/541269a
3. Open Science Collaboration. Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015). doi:10.1126/science.aac4716
4. Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* **533**, 452–454 (2016). (2016). doi:10.1038/533452a
5. Aschwanden, C. Science isn't broken. FiveThirtyEight. (2015). Available at: <https://fivethirtyeight.com/features/science-isnt-broken/>. (Accessed: 24th January 2017)
6. Buck, S. Solving reproducibility. *Science* **348**, 1403–1403 (2015). doi:10.1126/science.aac8041
7. Frew, J. & Dozier, J. Environmental informatics. *Annual Review of Environment and Resources* **37**, 449–472 (2012). doi:10.1146/annurev-environ-042711-121244
8. Jones, M. B., Schildhauer, M. P., Reichman, O. J. & Bowers, S. The new bioinformatics: Integrating ecological data from the gene to the biosphere. *Annual Review of Ecology, Evolution, and Systematics* **37**, 519–544 (2006).
9. Michener, W. K. & Jones, M. B. Ecoinformatics: Supporting ecology as a data-intensive science. *Trends in Ecology & Evolution* **27**, 85–93 (2012). doi:10.1016/j.tree.2011.11.016
10. Check Hayden, E. Mozilla plan seeks to debug scientific code. *Nature News* **501**, 472 (2013). doi:10.1038/501472a
11. Boettiger, C., Chamberlain, S., Hart, E. & Ram, K. Building software, building community: Lessons from the rOpenSci project. *Journal of Open Research Software* **3**, (2015). doi:10.5334/jors.bu
12. Wilson, G. *et al.* Good enough practices in scientific computing. *arXiv:1609.00037 [cs]* (2016).
13. Wilson, G. V. Where's the real bottleneck in scientific computing? *Sigma Xi, The Scientific Research Society* **Jan-Feb**, 5–6 (2006).
14. Baker, M. Scientific computing: Code alert. *Nature* **541**, 563–565 (2017). doi:10.1038/nj7638-563a

- 462 15. Barone, L., Williams, J. & Micklos, D. Unmet needs for analyzing biological big  
463 data: A survey of 704 NSF principal investigators. *bioRxiv* 108555 (2017).  
464 doi:[10.1101/108555](https://doi.org/10.1101/108555)
- 465 16. Wolkovich, E. M., Regetz, J. & O'Connor, M. I. Advances in global change research  
466 require open science by individual researchers. *Global Change Biology* **18**, 2102–  
467 2110 (2012). doi:10.1111/j.1365-2486.2012.02693.x
- 468 17. Nosek, B. A. *et al.* Promoting an open research culture. *Science* **348**, 1422–1425  
469 (2015). doi:10.1126/science.aab2374
- 470 18. Reichman, O. J., Jones, M. B. & Schildhauer, M. P. Challenges and opportunities of  
471 open data in ecology. *Science* **331**, 703–705 (2011). doi:10.1126/science.1197962
- 472 19. Shade, A. & Teal, T. K. Computing workflows for biologists: A roadmap. *PLOS Biol*  
473 **13**, e1002303 (2015). doi:10.1371/journal.pbio.1002303
- 474 20. Goodman, A. *et al.* Ten simple rules for the care and feeding of scientific data.  
475 *PLOS Computational Biology* **10**, e1003542 (2014).  
476 doi:10.1371/journal.pcbi.1003542
- 477 21. Sandve, G. K., Nekrutenko, A., Taylor, J. & Hovig, E. Ten simple rules for  
478 reproducible computational research. *PLOS Computational Biology* **9**, e1003285  
479 (2013). doi:10.1371/journal.pcbi.1003285
- 480 22. White, E. P. *et al.* Nine simple ways to make it easier to (re)use your data. *Ideas in*  
481 *Ecology and Evolution* **6**, (2013).
- 482 23. Kervin, K., Michener, W. & Cook, R. Common errors in ecological data sharing.  
483 *Journal of eScience Librarianship* **2**, (2013). doi:10.7191/jeslib.2013.1024
- 484 24. Lewandowsky, S. & Bishop, D. Research integrity: Don't let transparency damage  
485 science. *Nature News* **529**, 459 (2016). doi:10.1038/529459a
- 486 25. Michener, W. K. Ten simple rules for creating a good data management plan.  
487 *PLOS Computational Biology* **11**, e1004525 (2015).  
488 doi:10.1371/journal.pcbi.1004525
- 489 26. Mislan, K. A. S., Heer, J. M. & White, E. P. Elevating the status of code in ecology.  
490 *Trends in Ecology & Evolution* **31**, 4–7 (2016). doi:10.1016/j.tree.2015.11.006
- 491 27. Kratz, J. & Strasser, C. Data publication consensus and controversies.  
492 *F1000Research* (2014). doi:[10.12688/f1000research.3979.3](https://doi.org/10.12688/f1000research.3979.3)
- 493 28. Munafò, M. R. *et al.* A manifesto for reproducible science. *Nature Human*  
494 *Behaviour* **1**, 0021 (2017). doi:10.1038/s41562-016-0021
- 495 29. Martinez, C. *et al.* *Reproducibility in science: A guide to enhancing reproducibility*  
496 *in scientific results and writing*. (2014). Available at:  
497 <http://ropensci.github.io/reproducibility-guide/>

- 498 30. Tuyl, S. V. & Whitmire, A. L. Water, water, everywhere: Defining and assessing  
499 data sharing in academia. *PLOS ONE* **11**, e0147942 (2016).  
500 doi:10.1371/journal.pone.0147942
- 501 31. Baker, M. Why scientists must share their research code. *Nature* (2016).  
502 doi:10.1038/nature.2016.20504
- 503 32. Kidwell, M. C. *et al.* Badges to acknowledge open practices: A simple, low-cost,  
504 effective method for increasing transparency. *PLOS Biology* **14**, e1002456 (2016).  
505 doi:10.1371/journal.pbio.1002456
- 506 33. Wickham, H. & Grolemund, G. *R for data science*. (O'Reilly, 2016). Available at:  
507 <http://r4ds.had.co.nz/>
- 508 34. Lowndes, J. S. S. *et al.* Best practices for assessing ocean health in multiple  
509 contexts using tailorable frameworks. *PeerJ* **3**, e1503 (2015).  
510 doi:10.7717/peerj.1503
- 511 35. Lowndes, J. S. S. A biography of the ocean health index. OHI-science. (2017).  
512 Available at: <http://ohi-science.org/news/Biography-OHI>. (Accessed: 26th January  
513 2017)
- 514 36. Halpern, B. S. *et al.* An index to assess the health and benefits of the global ocean.  
515 *Nature* (2012). doi:10.1038/nature11397
- 516 37. Halpern, B. S. *et al.* Patterns and emerging trends in global ocean health. *PLoS*  
517 *ONE* **10**, e0117863 (2015). doi:10.1371/journal.pone.0117863
- 518 38. Ocean Health Index Team. Five years of global ocean health index assessments.  
519 (2016). Available at: <http://ohi-science.org/ohi-global/>.
- 520 39. Goldfuss, C. & Holdren. The nation's first ocean plans. The white house of  
521 president barack obama. (2016). Available at:  
522 [https://obamawhitehouse.archives.gov/blog/2016/12/07/nations-first-ocean-](https://obamawhitehouse.archives.gov/blog/2016/12/07/nations-first-ocean-plans)  
523 [plans](https://obamawhitehouse.archives.gov/blog/2016/12/07/nations-first-ocean-plans). (Accessed: 24th January 2017)
- 524 40. Hampton, S. E. *et al.* The tao of open science for ecology. *Ecosphere* **6**, 1–13  
525 (2015). doi:10.1890/ES14-00402.1
- 526 41. Mills, B. Introducing mozilla science study groups. Mozilla science. (2015).  
527 Available at: [https://science.mozilla.org/blog/introducing-mozilla-science-study-](https://science.mozilla.org/blog/introducing-mozilla-science-study-groups)  
528 [groups](https://science.mozilla.org/blog/introducing-mozilla-science-study-groups). (Accessed: 2nd August 2016)
- 529 42. R Core Team. *R: A language and environment for statistical computing*. R  
530 Foundation for Statistical Computing (2016). <https://www.R-project.org/>
- 531 43. RStudio Team. *RStudio: Integrated development for R*. RStudio, Inc, (2016).  
532 [www.rstudio.com](http://www.rstudio.com)
- 533 44. Git Team. *Git version control system*. (Git, 2016). <https://git-scm.com/>



534 45. GitHub. *GitHub: A collaborative online platform to build software*. (GitHub, 2016).  
535 <https://github.com>

536 46. Wilson, G. V. Software carpentry: Getting scientists to write better code by  
537 making them more productive. *Computing in Science and Engineering* 66–69 (2006).

538 47. Broman, K. Initial steps toward reproducible research. (2016). Available at:  
539 <http://kbroman.org/steps2rr/>.

540 48. McKiernan, E. C. *et al.* How open science helps researchers succeed. *eLife* **5**,  
541 e16800 (2016). doi:10.7554/eLife.16800

542 49. Seltenrich, N. Scaling the heights of data science. *Breakthroughs: The magazine of*  
543 *the UC Berkeley College of Natural Resources* **Fall**, (2016). Available at:  
544 <https://nature.berkeley.edu/breakthroughs/opensci-data>

545 50. Wilson, G. *et al.* Best practices for scientific computing. *PLOS Biol* **12**, e1001745  
546 (2014). doi:10.1371/journal.pbio.1001745

547 51. Haddock, S. H. & Dunn, C. W. *Practical computing for biologists*. (Sinauer  
548 Associates, Inc., 2011). <http://practicalcomputing.org/>

549 52. Barnes, N. Publish your computer code: It is good enough. *Nature* **467**, 753  
550 (2010).

551 53. Frazier, M., Longo, C. & Halpern, B. S. Mapping uncertainty due to missing data in  
552 the global ocean health index. *PLOS ONE* **11**, e0160377 (2016).  
553 doi:10.1371/journal.pone.0160377

554 54. ESRI. *ArcGIS platform*. (2016). <http://www.esri.com/software/arcgis>

555 55. The QGIS Team. *QGIS project*. (2016). <http://www.qgis.org>

556 56. The Python Team. *Python.org*. (2016).

557 57. Wickham, H. Tidy data. *Journal of Statistical Software* **59**, 1–23 (2014).  
558 doi:10.18637/jss.v059.i10

559 58. Wickham, H. Tidyverse tidyweb. (2017). Available at: <http://tidyverse.org/>.  
560 (Accessed: 29th January 2017)

561 59. Wickham, H. *Tidyverse: Easily install and load 'tidyverse' packages*. (2016).  
562 Available at: <https://CRAN.R-project.org/package=tidyverse>

563 60. Fischetti, T. How dplyr replaced my most common r idioms. StatsBlogs. (2014).  
564 Available at: [http://www.onthelambda.com/2014/02/10/how-dplyr-replaced-my-](http://www.onthelambda.com/2014/02/10/how-dplyr-replaced-my-most-common-r-idioms/)  
565 [most-common-r-idioms/](http://www.onthelambda.com/2014/02/10/how-dplyr-replaced-my-most-common-r-idioms/). (Accessed: 22nd September 2016)

566 61. RStudio Team. R Markdown. (2016). Available at:  
567 <http://rmarkdown.rstudio.com/>.

568 62. Allaire, J.J. *et al.* *R Markdown: Dynamic documents for R*. (2016). Available at:  
569 <https://CRAN.R-project.org/package=rmarkdown>

570 63. Ocean Health Index Team. *Ocean Health Index ohicore package*. (2016).

571 64. Wickham, H. *R packages*. (O'Reilly, 2015). Available at: <http://r-pkgs.had.co.nz/>

572 65. Wickham, H. & Chang, W. *Devtools: Tools to make developing R packages easier*.  
573 (2016). Available at: <https://CRAN.R-project.org/package=devtools>

574 66. Wickham, H., Danenberg, P. & Eugster, M. *Roxygen2: In-source documentation for*  
575 *R*. (2015). Available at: <https://CRAN.R-project.org/package=roxygen2>

576 67. Ram, K. Git can facilitate greater reproducibility and increased transparency in  
577 science. *Source Code for Biology and Medicine* **8**, 7 (2013). doi:10.1186/1751-0473-  
578 8-7

579 68. Blischak, J. D., Davenport, E. R. & Wilson, G. A quick introduction to version  
580 control with git and GitHub. *PLOS Comput Biol* **12**, e1004668 (2016).  
581 doi:10.1371/journal.pcbi.1004668

582 69. Perez-Riverol, Y. *et al.* Ten simple rules for taking advantage of git and GitHub.  
583 *PLOS Comput Biol* **12**, e1004947 (2016). doi:10.1371/journal.pcbi.1004947

584 70. Duhigg, C. What google learned from its quest to build the perfect team. *The New*  
585 *York Times* (2016). Available at:  
586 [http://www.nytimes.com/2016/02/28/magazine/what-google-learned-from-its-](http://www.nytimes.com/2016/02/28/magazine/what-google-learned-from-its-quest-to-build-the-perfect-team.html)  
587 [quest-to-build-the-perfect-team.html](http://www.nytimes.com/2016/02/28/magazine/what-google-learned-from-its-quest-to-build-the-perfect-team.html)

588 71. Perkel, J. Democratic databases: Science on GitHub. *Nature* **538**, 127–128  
589 (2016). doi:10.1038/514127a

590 72. Casadevall, A. & Fang, F. C. Reproducible science. *Infection and Immunity* **78**,  
591 4972–4975 (2010).

592 73. Wilson, G. Software carpentry: Lessons learned. *F1000Research* (2016).  
593 doi:[10.12688/f1000research.3-62.v2](https://doi.org/10.12688/f1000research.3-62.v2)

594 74. Hampton, S. E. *et al.* Big data and the future of ecology. *Frontiers in Ecology and*  
595 *the Environment* **11**, 156–162 (2013). doi:10.1890/120103

596 75. Lohr, S. For big-data scientists, ‘janitor work’ is key hurdle to insights. *The New*  
597 *York Times* (2014). Available at:  
598 [http://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-](http://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html)  
599 [to-insights-is-janitor-work.html](http://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html)

600 76. FitzJohn, R., Pennell, M., Zanne, A. & Cornell, W. Reproducible research is still a  
601 challenge. *ROpenSci*. (2014). Available at: [/blog/2014/06/09/reproducibility](http://blog/2014/06/09/reproducibility).  
602 (Accessed: 2nd August 2016)

- 603 77. Boland, M. R., Karczewski, K. J. & Tatonetti, N. P. Ten simple rules to enable  
604 multi-site collaborations through data sharing. *PLOS Computational Biology* **13**,  
605 e1005278 (2017). doi:10.1371/journal.pcbi.1005278
- 606 78. Perkel, J. M. Scientific writing: The online cooperative. *Nature* **514**, 127–128  
607 (2014). doi:10.1038/514127a
- 608 79. Rbloggers. How twitter improved my ecological model. Rbloggers. (2015).  
609 Available at: [https://www.r-bloggers.com/how-twitter-improved-my-ecological-](https://www.r-bloggers.com/how-twitter-improved-my-ecological-model/)  
610 [model/](https://www.r-bloggers.com/how-twitter-improved-my-ecological-model/). (Accessed: 26th January 2017)

611

## Acknowledgements

The Ocean Health Index is a collaboration between Conservation International and the National Center for Ecological Analysis and Synthesis at the University of California at Santa Barbara. We thank Johanna Polsenberg, Steve Katona, Erich Pacheco, and Lindsay Mosher who are our partners at Conservation International. We thank all past contributors and funders that have supported the Ocean Health Index, including Beau and Heather Wrigley and The Pacific Life Foundation. We also thank all the individuals and groups that openly make their data, tools, and tutorials freely available to others. Finally, we thank Hadley Wickham, Karthik Ram, Kara Woo, and Mark Schildhauer for friendly review of the developing manuscript.

See [ohi-science.org/betterscienceinlesstime](https://ohi-science.org/betterscienceinlesstime) as an example of a website built with RMarkdown and the RStudio-GitHub workflow, and for links and resources referenced in the paper.

## Author Contributions

All authors developed concepts and wrote the paper.

## Competing interests

The authors declare no competing financial interests.

## Corresponding author

Correspondence to Julia Lowndes: [lowndes@nceas.ucsb.edu](mailto:lowndes@nceas.ucsb.edu)

## Figures

### Figure 1

**Better science in less time, illustrated by the Ocean Health Index project.** Every year since 2012 we have repeated Ocean Health Index (OHI) methods to track change in global ocean health<sup>35,36</sup>. Increased reproducibility and collaboration has reduced the amount of time required to repeat methods (size of bubbles) with updated data annually, allowing us to focus on improving methods each year (biggest innovations written as text). The original assessment in 2012 focused solely on scientific methods (e.g., obtaining and analyzing data; developing models; calculating and presenting results (dark shading)). In 2013, by necessity we gave more focus to data science (e.g., data organisation and wrangling; coding; versioning; documentation (light shading)), using open data science tools. We established R as the main language for all data preparation and modeling (using RStudio), which drastically decreased the time involved to complete the assessment. In 2014, we adopted Git and GitHub for version control, project management, and collaboration. This further decreased the time required to repeat the assessment. We also created the OHI Toolbox, which includes our R package ohicore for core analytical operations used in all OHI assessments. In subsequent years we have continued (and plan to continue) this trajectory towards better science in less time by improving code with principles of tidy data<sup>33</sup>; standardising file and data structure; and focusing more on communication, in part by creating websites with the same open data science tools and workflow. See text and Table 1 for more details.

# Tables

**Table 1**

## Summary of the primary open data science tools used to upgrade

**reproducibility, collaboration, and communication, by task.** The transition to

using open data science tools was incremental (see Figure 1). All tasks are

accomplished with the RStudio–GitHub workflow that is underpinned by R and Git.

This workflow streamlines collaboration by capturing each individual’s contribution

to the project – thus taking care of bookkeeping – for tasks from data processing and

analysis to creating documents and websites with embedded results that are

updatable. Note that collaboration is not only for labs and teams, but also for each

individual’s ‘future self’.

	Task	Then	Now	Primary open data science tools
<b>Reproducibility</b>	data preparation	manually (i.e., Excel)	coded in R	R packages: tidyverse (dplyr, tidyr, ggplot2). Documentation: R Markdown
	modeling	multiple programming languages	R functions and ohicore package	R packages: tidyverse, devtools, roxygen2, git2r
	version control	file duplication and renaming	Git	Git; interface with Git and GitHub primarily through RStudio
	organisation	individual conventions	standardised team convention	RStudio projects, GitHub repositories. File structure protocols
	coding	separate languages and conventions	R; standardised team convention	Principles of tidy data; tidyverse
<b>Collaboration</b>	workflow and project management	individual conventions	(simplified) GitHub workflow	GitHub, RStudio
	internal	email	centralised,	GitHub issues

	collaboration		archived conversations	
	sharing data	ftp download	all versions and Releases available online	ohi-science.org/ohi- global
<b>Communication</b>	sharing methods	published manuscript and supplementary material	published on our website	ohi-science.org website, with linked R Markdown outputs (webpages, presentations, etc)

666

## 667 **Boxes**

### 668 **Box 1**

669 **Resources to learn open data science tools.** These are some of the free, online  
670 resources that we used to learn and develop a workflow with R, RStudio, Git, and  
671 GitHub. These resources exposed us to what was possible, and helped us build skills  
672 to incorporate concepts and tools into our own workflow. This is by no means an  
673 exhaustive list. See also Box 2 for strategies of how to get started.

674

#### 675 **Primarily R**

676 **R for Data Science** book by Hadley Wickham and Garrett Grolemund  
677 ([r4ds.had.co.nz](http://r4ds.had.co.nz))

678 **RStudio's webinars** on-demand videos by RStudio  
679 ([rstudio.com/resources/webinars](http://rstudio.com/resources/webinars))

680 **RStudio's cheatsheets** PDFs by RStudio  
681 ([rstudio.com/resources/cheatsheets](http://rstudio.com/resources/cheatsheets))

682 **CRAN Task Views** to identify useful packages by category of task  
683 ([cran.r-project.org/web/views](http://cran.r-project.org/web/views))

684 **R Packages** book by Hadley Wickham  
685 ([r-pkgs.had.co.nz](http://r-pkgs.had.co.nz))

#### 686 **Combination RStudio-GitHub**

687 **Happy Git With R** short-course by Jenny Bryan  
688 ([happygitwithr.com](http://happygitwithr.com))

689 **UBC Stats545: Data Wrangling, Exploration, and Analysis with R** university  
690 course by Jenny Bryan  
691 ([stat545-ubc.com](http://stat545-ubc.com))

692 **Software Carpentry** workshops, teaching and learning communities  
693 ([software-carpentry.org](http://software-carpentry.org))

694 example 2-day course: "Reproducible Science with RStudio and GitHub"  
695 [jules32.github.io/2016-07-12-Oxford/overview](https://jules32.github.io/2016-07-12-Oxford/overview)

#### 696 **Community discussion**

697 **#rstats on Twitter** online discussions  
698 ([twitter.com/search?q=%23rstats&src=typd](https://twitter.com/search?q=%23rstats&src=typd))

699 **Not So Standard Deviations** podcast by Roger Peng and Hilary Parker  
700 ([soundcloud.com/nssd-podcast](http://soundcloud.com/nssd-podcast))

701 **R-Bloggers** blog  
702 ([r-bloggers.com](http://r-bloggers.com))



703 **RStudio** blog  
704 ([blog.rstudio.org](https://blog.rstudio.org))  
705 **Data Carpentry** blog  
706 ([datacarpentry.org/blog](https://datacarpentry.org/blog))  
707

## Box 2

**Strategies to learn in an intentional way.** The resources listed in Box 1 have helped us learn open data science principles and tools in an intentional way: We felt empowered (vs. panicked), we learned to think ahead (vs. quick fixes for single purposes), and we learned with a community (vs. in isolation). There is a whole ecosystem of open data science principles, practices, and tools (including R, RStudio, Git, and GitHub) and no single way to begin learning. These are a few strategies you can consider as you get engaged.

### **Self-paced learning**

Box 1 lists resources to learn open data science principles and tools that you can use at your own pace. The books and courses provide in-depth philosophies and are good for initial learning as well as for reference later on. Webinars and podcasts are generally under an hour.

### **Join and/or create communities**

Learning together and supporting each other peer-to-peer can be more fun and rewarding. You can become a "champion" for others by showing leadership as you learn. Start off by watching a webinar with a friend or group during lunch or a happy hour. Learn enough about a useful R package to share in your lab meetings; you learn best by teaching. In traditional journal clubs or lab meetings, discuss an academic article on importance of reproducibility, collaboration, and coding<sup>14,22,69,78</sup>. Search if your institution or city has local Meetup.com groups, or create your own.

Additionally, join or keep tabs on communities online. Mozilla Study Groups are a network of 'journal-clubs' where scientists teach scientists computing skills ([science.mozilla.org/programs/studygroups/join](https://science.mozilla.org/programs/studygroups/join)). rOpenSci is a developer collective building R-based tools to facilitate open science ([ropensci.org](https://ropensci.org)). Also look on Twitter for #rstats discussions and then follow individuals from those conversations.

### **Ask for help**

Local and online communities are a great resource to ask when you need help. Expecting that someone has already asked your question can help you both articulate the problem clearly and identify useful answers. Often, pasting error messages directly into Google will get you to the best answers quickly. Many answers come from online forums, including

742 [StackOverflow.com](#)<sup>14</sup>, or even Twitter itself (e.g., 'How Twitter Improved My  
743 Ecological Model')<sup>79</sup>.

#### 744 **Attend in-person workshops and conferences**

745 In-person workshops can be extremely valuable and give you an opportunity  
746 to get direct help from instructors and helpers. Software Carpentry and Data  
747 Carpentry run 2-day bootcamps that teach skills for research computing; you  
748 can attend a scheduled workshop or request your own ([software-](#)  
749 [carpentry.org](#); [datacarpentry.org](#)). Attend conferences like useR (example:  
750 [user2017.brussels](#)) both for skill-building and to learn how others are using  
751 these tools.

#### 752 **Watch presentations from past conferences**

753 More and more, slide decks and videos of presentations are appearing online.  
754 For example, you can see presentations from the the 2016 useR conference  
755 ([user2016.org](#)) and the 2017 RStudio conference ([rstudio.com/conference](#)).

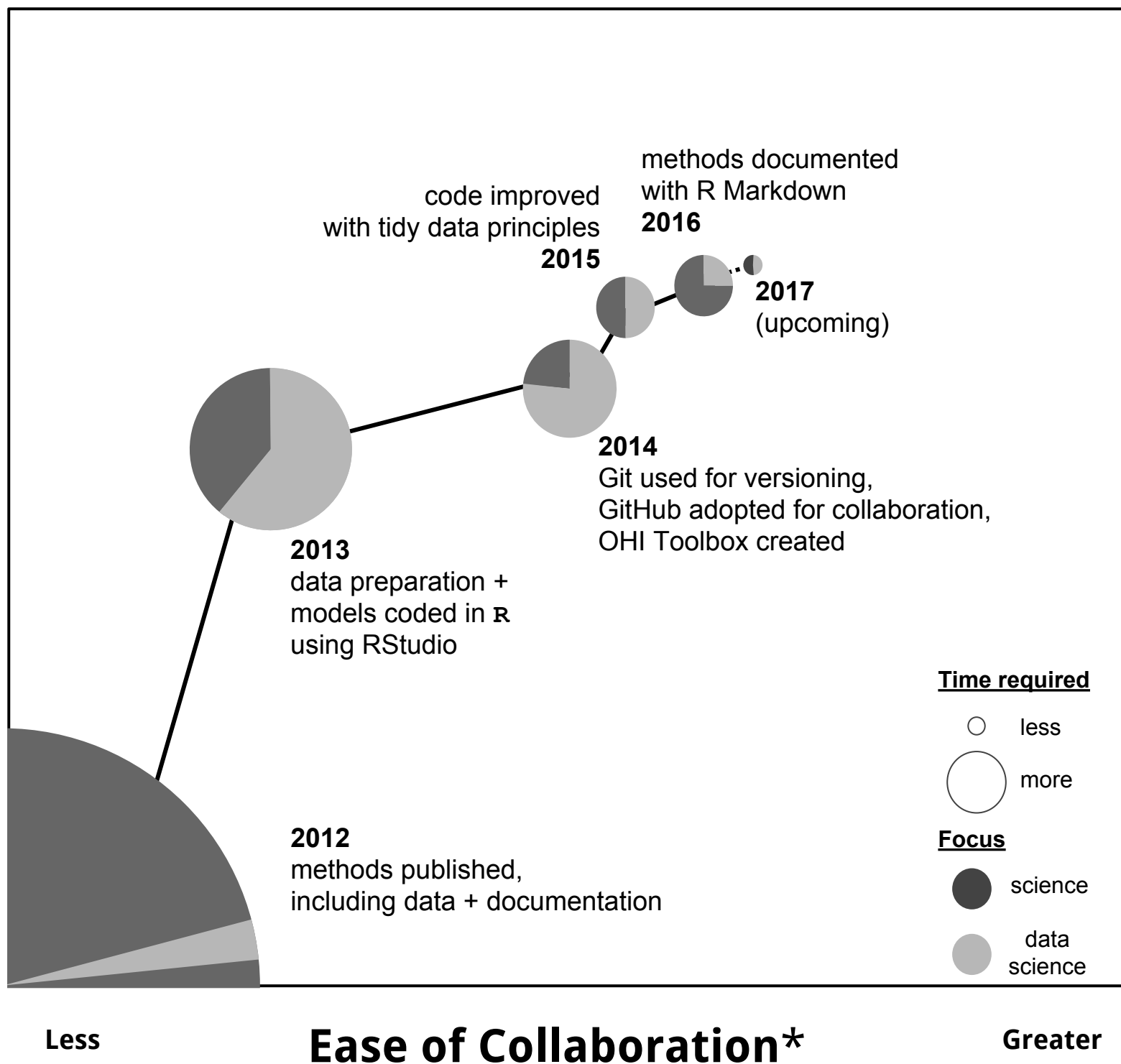
#### 756 **Read Blogs**

757 There are many individuals who blog about open data science concepts, R  
758 packages, workflows, etc. Try Googling a package you're using, or going to  
759 the website of someone you are following on Twitter.

Reasonably  
easily

Ease of Reproducibility

With great  
difficulty



\*including future self