# Moving toward a sustainable ecological science: don't let data go to waste!

Timothée Poisot, Dominique Gravel

Nov. 2012

## Intro

Claude Bernard (Bernard 1864) wrote that "art is *me*; science is *us*". This sentence has two meaning. First, the altruism of scientists is worth more than the self-indulgence of mid-nineteenth century Parisian art scene. Second, and we will keep this one in mind, creativity and insights come from individuals, but validation and rigor are reached through collective efforts, cross- validation, and peerage. Given enough time, the conclusions reached and validated by the efforts of many will take prominence over individualities, and this (as far as Bernard is concerned), is what science is about. With the technology available to a modern scientist, one should expect that the dissolution of *me* would be accelerated, and that several scientists should be able to cast a critical eye on data, and use this collective effort to draw robust conclusions.

In molecular evolution, there exists a large number of databases (GenBank, EMBL, SwissProt, and many more) in which information can be retrieved. This values (and allows) a new type of scientific research: building over the raw material of others, it is now possible to identify new phenomenon or evaluate the generality of previously studied ones. The job of these scientists is not to *make* data, neither to *stole* them, it's rather to gather them and, most of all, look at them in a different way. This would not be possible, if not for the existence of public, free, online repositories. It's impossible to be as enthusiastic when looking at current practices in ecology. Apart from a few, non-specific initiatives (*DataDryad*), or small-scale initiatives which are not always properly maintained (*Interaction Web Networks Database*), there is no data sharing culture among ecologists.

In this paper, using the example of ecological networks, we will argue that improving our data sharing practices will improve both the science, and the reputation of the scientists. We will illustrate how simple steps can be taken to greatly improve the situation, and how we can encourage the practice of data-sharing at different levels.

## Why we morally must

- The data are publicly funded and do not belong to the researcher who collected them

- It allows reproducibility of the science, which is supposed to be the rule

Using journals to publish scientific information should not only serve the purpose of disseminating an interesting discussion of data: it should maximize the ability of other researchers to replicate, and thus both validate and expand, results. It's interesting to see that, while editors and referees alike are very careful about the way the *Materials & Methods* parts of a paper are worded, it's extremely rare to recieve any comment about the data availability. This can cause problems at all steps of the life of a paper. How can a paper describing a new method be adequately reviewed if data are not available? How can you be sure that you are correctly applying a method if you can't reproduce the results? The movement of *reproducible research* advocates that a paper can be self-contained, *i.e.* be not only the text, but also the data, and the computer code to reproduce the figures. Even without going to such lengths, releasing data and computer code alongside a paper should be viewed as an ethical decision. Barnes [@Barnes2010] made the point that computer code is good enough to be shared, even though researchers are not professional programmers.

- It will fight bad authorship practices, people hitch-hicking on other people's work

- Data are costly (time and money) to acquire, acquiring new instead of using old ones is wasteful

(Heidorn 2008) dark data, there is already enough material to answer some pending questions

## Why is it beneficial for the one who collected data

- A proxy to your science: data are a mean for people to get familiar with what you do

[@Ince2012] improves reproducibility and adequate communication of your results

(Vandewalle 2012) showed that sharing computer code improved the scientific impact

- It stimulates collaboration and creativity

- A measure of your productivity

## How we technically can

- End the rule of Excel: JSON schemes or XML to represent context-rich data

- FigShare and other projects: data can have a DOI and be cited/shared

- Local databases but linked globally: APIs and programmatic access

## How it should be encouraged

- Journal policies, and referee expertise

Several journals are now asking the authors to deposit the data in a public repository. This is mandatory for sequences (GenBank), and various journals recommend the use of TreeBase, DataDryad, or FigShare. The referees are, however, rarely asked to evaluate if the adequate data are released (*e.g.* network metrics and summary statistics instead of full networks), and even more rarely given access to the data during the evaluation process.

- Evaluation for funding

## Conclusion

## List of possible boxes

- The story of the BCI data

- What we could tell about network biogeography with public data?

Bernard, C. 1864. "Introduction à l'étude de la médecine expérimentale." *Revue des cours scientifiques. Paris.* http://www.bouquineux.com/pdf/Bernard-Introduction\_a\_l\_etude\_de\_la\_medecine\_experimentale.pdf.

Heidorn, P. Bryan. 2008. "Shedding Light on the Dark Data in the Long Tail of Science." *Library Trends* 57: 280–299. doi:10.1353/lib.0.0036. http://muse.jhu.edu/content/crossref/journals/library\_trends/v057/57.2.heidorn.html.

Vandewalle, Patrick. 2012. "Code Sharing is Associated with Research Impact in Image Processing." *Computing in Science and Engineering*: 1–5. http://ieeexplore.ieee.org/xpls/abs\_all.jsp?arnumber=6200247.