

1 Building bridges between data providers and users: best practices and lessons  
2 learned

3 foo bar<sup>\*,a</sup>

4 <sup>a</sup>*rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA*

5 **Abstract**

6 Corresponding Author:

7 foo bar

8 rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA

9 Email address: [stuff@ropensci.org](mailto:stuff@ropensci.org)

---

\*Corresponding author  
Email address: `stuff(at)ropensci.org` (foo bar)

abstract text ...

## 10 Introduction

11 intro text ...

## 12 OUTLINE (from chat btw Carl/Scott)

- 13 • use cases necessitate FTP vs. csv dump vs. rest API vs etc.
- 14 • within APIS: what are best practices
- 15 • discoverability - how do you find out if an API has data (taxonomic, geospatial, etc.) you want?
- 16 • data formats: tabular vs. JSON/XML, etc.
- 17 • combining data: e.g., using identifiers instead of names
- 18 • best practices in relational databases, briefly
- 19 • ropensci filling gaps btwn data providers and scientists
  - 20 – communicating from data providers to sci.
  - 21 – communicating from sci. to data providers

## 22 Overview of the landscape

23 There is an increasing amount of data available to researchers. Leveraging that data efficiently and  
24 reproducibly ideally requires software. However, researchers have limited time - thus, most rightly focus  
25 on the science. Furthermore, we want researchers to focus on science. Given this, there are a number of  
26 issues that others must handle:

- 27 • We need well made open source software to help researchers leverage data
- 28 • We need people as bridges between data providers and data users (researchers)
- 29 • Slack needs to be taken up to help small scientific databases
- 30 • We need best practices for working with data

## 31 Matching the data format to the use case

32 There are innumerable data formats, and there is no one data format that is best for every situation.  
33 Ideally, one leverages the right data format for their use case. The following use cases highlight some of  
34 the diversity and what data formats they best match.

- 35 • Large amount of data: This varies surely, and depends on connection speeds, computers available,  
36 etc. - but for sake of argument lets say that  $> 1$  GB is large data for this use case. It doesn't  
37 make a lot of sense to provide this through an API, and makes sense to provide as a compressed  
38 format. Data of this type makes sense to provide via FTP or similar (Amazon S3, http file server,  
39 etc.).
- 40 • Small amount of data: Probably the majority of data use cases are "small data". For example, a  
41 spreadsheet with 100 or 100,000 rows is small data. The delivery mechanism can be more flexible  
42 with this kind of data. You can definitely serve this data over FTP, but can also simply provide  
43 csv/tsv files, and can serve the data over an API.
- 44 • Data constantly changing: This is a good use case for delivering data via an API. APIs connect  
45 to an underlying database that can change as much as the data providers need. However, the API  
46 ideally changes only very slowly so that clients can depend on the interface. It's easier to update  
47 data incrementally over an API than if there's small changes in lots of files on an FTP server.

## 48 Discoverability

49 Discovering what kind of data you can get from a data source before actually getting that data is an  
50 important one. For example, say there's a database with 1000 GB of data. You don't want to have to  
51 download that database to your own machine, then search through it to find the data you want. Ideally  
52 there's a fast way to query a database (perhaps it's metadata) before delving into the data fetching  
53 process that will likely take longer.

54 Unfortunately, most datasets are very lacking in metadata. However, when metadata is well filled out,  
55 it makes data discovery much easier. For example, ...

56 Suggestions?: do x, y, and z

57 **Bridges between data providers and users**

58 **Acknowledgments**

59 This project was supported in part by the Alfred P Sloan Foundation (Grant No. G-2014-13485), and  
60 in part by the Helmsley Foundation (Grant No. 2016PG-BRI004).

61 **References**