

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

3 foo bar^{*,a}

4 ^a*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

*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 that 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 provdiers to sci.
 - 21 – communicating from sci. to data provdiers

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 tacken up to help small scientific databases
- 30 • We need best practices for working with data

31 Data formats

32 Data formats are incredibly diverse. In addition, there's many that are used more often in scientific use
33 cases than elsewhere (e.g. NetCDF). Data formats lead to many problems, often in some only working
34 on some platforms, and in translating between data formats.

- 35 • Tabular (csv/tsv)
- 36 • JSON
- 37 • XML
- 38 • PDF

39 Matching the data format to the use case

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

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

56 Discoverability

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

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

64 Suggestions?: do x, y, and z

65 Bridges between data providers and users

66 Data providers - sysadmin's, software engineers, data curators, publishers, domain experts, and more
67 - are keenly focused on providing a great product for researchers and the public. Data providers in
68 the scientific space are not usually very profitable. Thus, time is very limited. In addition, the very
69 technical bent of the data providers may not be the best fit for researchers that don't understand the
70 same terminology, etc.

71 This leaves space for bridges to be built between data providers and researchers - ideally people that
72 understand both data providers and researchers and can speak both languages. rOpenSci has been
73 doing this for a number of communities.

74 XXXX

75 Acknowledgments

76 This project was supported by the Helmsley Foundation (Grant No. 2016PG-BRI004).

77 References