

GigaScience

Keeping it light: (Re)analyzing community-wide datasets without major infrastructure --Manuscript Draft--

| | | |
|---|--|-------------------|
| Manuscript Number: | GIGA-D-18-00196R2 | |
| Full Title: | Keeping it light: (Re)analyzing community-wide datasets without major infrastructure | |
| Article Type: | Commentary | |
| Funding Information: | Gordon and Betty Moore Foundation (GBMF4551) | Dr C. Titus Brown |
| Abstract: | <p>DNA sequencing technology has revolutionized the field of biology, shifting biology from a data-limited to data-rich state. Central to the interpretation of sequencing data are the computational tools and approaches that convert raw data into biologically meaningful information. Both the tools and the generation of data are actively evolving, yet the practice of re-analysis of previously generated data with new tools is not commonplace. Re-analysis of existing data provides an affordable means of generating new information and will likely become more routine within biology, yet necessitates a new set of considerations for best practices and resource development. Here, we discuss several practices that we believe to be broadly applicable when re-analyzing data, especially when done by small research groups.</p> | |
| Corresponding Author: | Harriet Alexander, PhD UNITED STATES | |
| Corresponding Author Secondary Information: | | |
| Corresponding Author's Institution: | | |
| Corresponding Author's Secondary Institution: | | |
| First Author: | Harriet Alexander, PhD | |
| First Author Secondary Information: | | |
| Order of Authors: | Harriet Alexander, PhD Lisa K Johnson, BA C. Titus Brown, PhD | |
| Order of Authors Secondary Information: | | |
| Response to Reviewers: | We have corrected the Johnson et al citation to be in press at GigaScience. We have also removed one citation (Kodama 2012) to conform to the 10 citation limit. | |
| Additional Information: | | |
| Question | Response | |
| Are you submitting this manuscript to a special series or article collection? | No | |
| Experimental design and statistics | No | |
| Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available | | |

| | |
|---|--|
| <p>in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p> | |
| <p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p> <p>"</p> | <p>This is a</p> |
| <p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p> | <p>No</p> |
| <p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Resources</p> | <p>This is a commentary piece to accompany the piece recently submitted by Johnson et al. As such, there are no data being analyzed.</p> |

| | |
|---|------------|
| <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p> <p>"</p> | |
| <p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p> | <p>Yes</p> |

[Click here to view linked References](#)*GigaScience*, 2018, 1–3doi: [xx.xxxxx/xxxxx](#)Manuscript in Preparation
Commentary

COMMENTARY

Keeping it light: (Re)analyzing community-wide datasets without major infrastructure

Harriet Alexander^{1,2}, Lisa K. Johnson^{1,3} and C. Titus Brown^{1,3,4*}

¹Population Health and Reproduction, University of California, Davis, CA, USA and ²Biology Department, Woods Hole Oceanographic Institution, Woods Hole, MA, USA and ³Molecular, Cellular, and Integrative Physiology Graduate Group, University of California, Davis, CA, USA and ⁴Genome Center, University of California Davis, CA, USA

*ctbrown@ucdavis.edu

Abstract

DNA sequencing technology has revolutionized the field of biology, shifting biology from a data-limited to data-rich state. Central to the interpretation of sequencing data are the computational tools and approaches that convert raw data into biologically meaningful information. Both the tools and the generation of data are actively evolving, yet the practice of re-analysis of previously generated data with new tools is not commonplace. Re-analysis of existing data provides an affordable means of generating new information and will likely become more routine within biology, yet necessitates a new set of considerations for best practices and resource development. Here, we discuss several practices that we believe to be broadly applicable when re-analyzing data, especially when done by small research groups.

Key words: reproducibility; data reuse; open data;

Background

Advances in high-throughput, next-generation sequencing technologies have catapulted biology into a new computational era. In fields of biology that leverage sequencing data, the primary limiting step in the earlier stages of biological inquiry has increasingly shifted away from data generation to data analysis. Concomitant with the increasing emphasis on the computational processing of these data is the advancement of the computational tools available for such analyses: new computational approaches for the analysis of these data are constantly being created, tested, and proved worthy of use. Yet, outside of computational lab groups, the practice of re-analysis of previously generated data with new tools and approaches is not commonplace. Such re-analysis has great utility and will become more routine within the life sciences, yet re-analysis necessitates a new set of considerations for best practices and resource development.

Our interest in the issues surrounding re-analysis was spurred by a large-scale sequencing project: the Marine Mi-

crobial Transcriptome Sequencing Project (MMETSP), which generated 678 transcriptomes, spanning 396 different strains of eukaryotic microbial eukaryotes isolated from marine settings [1]. This dataset is an invaluable resource within the oceanographic community [2, 1], as it exponentially expands the accessible genetic information base of marine protistan life. Moreover, the MMETSP has created a uniquely useful test dataset for computational biologists. The MMETSP dataset spans a large evolutionary history of organisms, and all of the 678 transcriptomes were prepared and sequenced in a consistent way [2]. The sequencing project, which was completed in 2014, was originally assembled by the National Center for Genome Resources using a custom pipeline that employed the best available computational tools at the time [3, 4].

Since the original MMETSP analysis, new tools and techniques for the assembly of *de novo* transcriptomes from RNAseq data have been described and preexisting tools have been improved upon [5]. Moreover, new annotation tools and databases have become available. The transcriptome assembly project described in Johnson et al. [6] was designed to create

Compiled on: November 14, 2018.

Draft manuscript prepared by the author.

1 a streamlined and reproducible assembly framework that not
2 only enables the re-analysis of these datasets, but creates a
3 framework to facilitate easy and rapid re-analyses in the fu-
4 ture.

5 These secondary data products of sequencing, such as an-
6 notated assemblies, should be viewed as hypotheses gener-
7 ated from the underlying biology, rather than some immutable
8 “truth”. As such these secondary data products can continue
9 to be improved as new tools are developed. For example, we
10 note that MacManes [7] described several limitations and chal-
11 lenges of current assembly technology and developed an im-
12 proved Oyster River Protocol, which we could use to generate
13 another, perhaps improved, MMETSP assembly.

14 Ultimately, such iterations on the original raw data have
15 the potential to improve upon the secondary data products –
16 the assembled transcriptomes and associated annotations that
17 are relied upon by the broader community for biological inquiry.
18 Through this process, we developed several practices that we
19 believe to be broadly applicable when re-analyzing data, espe-
20 cially when done by small research groups.

21 Main text

22 Storage of secondary data products

23 Funding agencies and academic journals now mandate the de-
24 position of raw data into digital repositories (e.g. NCBI Se-
25 quence Read Archive (SRA) and Gene Expression Omnibus, Eu-
26 ropean Nucleotide Archive). Thus, to date, the majority of the
27 sequence data that has been generated and published is openly
28 available online for reference and use in other studies. The
29 sharing and availability of raw data from high-throughput se-
30 quencing studies has been largely managed through the de-
31 velopment of archival services such as the SRA, which was
32 established as part of the International Nucleotide Sequence
33 Database Collaboration (INSDC)[8]. The SRA currently con-
34 tains more than 1.8e16 bases of information (~7e15 are open
35 access)¹. While a tremendous resource for biological inquiry,
36 a major problem remains in that raw sequencing data is not
37 the most directly useful form of sequencing data. Rather, bi-
38 ologists rely heavily upon the computationally generated sec-
39 ondary products of sequencing reads (e.g. assembled transcrip-
40 tomes or genomes, annotations, associated count-based data,
41 etc.). There is a dearth of these secondary products in central,
42 publicly accessible databases, such as the Transcriptome Shot-
43 gun Assembly (TSA) Sequence Database.

44 In fact, a substantial proportion of these data products
45 might be aptly categorized as “dark data,” as they are largely
46 undiscoverable and often archived independently in associa-
47 tion with a publication or on private servers. Even more
48 limiting, however, is that the guidelines for public databases
49 such as the TSA specifically state that “Assemblies from se-
50 quences not directly sequenced by the submitter” should not
51 be uploaded to the TSA, thereby excluding the potential for
52 reassembled datasets to be made available and directly linked
53 to preexisting BioProjects, BioSamples, TSAs, and SRA entries
54 (<https://www.ncbi.nlm.nih.gov/genbank/tsa/>).

55 From the perspective of our MMETSP re-analysis, we argue
56 the community needs more than a place to put the primary
57 and secondary data products associated with a single publica-
58 tion. Ideally, the results of each re-analysis would be deposited
59 in a discoverable location, but would have a coherent archival
60 procedure that is lab-independent, easily searchable, and “for-
61 ward discoverable” (i.e. when a new version of a data product

is released, old versions can point to the new version). More-
over, such an archival platform would ideally document the full
provenance of the secondary data product. Movement towards
this kind of data archival system are being made both with the
development of alternative scientific data publication models
(e.g. the Research Object[9]) as well as integration of metadata
models (such as the Resource Description Framework) onto ex-
isting scientific databases like the European Bioinformatics In-
stitute (EBI) [10], but policies surrounding secondary data prod-
ucts will need to change.

22 Directly linking secondary data products to prove- 23 nance of work-flow

24 In the absence of a community database specifically for the
25 type of secondary product that was produced in this analysis,
26 we opted to upload the assemblies, annotations, and counts
27 to Zenodo (<https://zenodo.org>), a scientific data repository
28 founded by CERN, which provided a DOI for the assemblies
29 (<https://doi.org/10.5281/zenodo.740440>). The header informa-
30 tion for each assembly was modified to contain the DOI. We
31 then created a Github repository containing the scripts used
32 to generate the assemblies. The repository was then archived
33 with Zenodo, which generated a single DOI for the project
34 (<https://doi.org/10.5281/zenodo.594854>).²

35 As such, the scripts used in the generation of transcrip-
36 tomes are directly linked through a unique DOI to the data prod-
37 ucts that are listed in the directory. Since the scripts are easily
38 accessible, they can be tweaked to re-analyze the primary se-
39 quence data using different parameters or tools, and the new
40 pipeline and output files can be archived again with Zenodo us-
41 ing the same approach as above. Moreover, the Zenodo archival
42 system will then automatically indicate the presence of other
43 versions of a given repository such that a user might be sure to
44 use the newest version of an assembly. In the future, such an
45 approach might be further complemented by the integration of
46 a JSON Linked Data file detailing the metadata for the assembly
47 product, such as the pipeline used and previous versions of the
48 assemblies.³

49 Conclusion

50 The Github-Zenodo framework presented here represents an
51 efficient way for small research groups (i.e. a graduate student)
52 to host and link both the code and results from large-scale
53 re-analysis projects in a publicly accessible way. The direct
54 linking of protocols and metadata to output data products is
55 paramount in the data heavy future of scientific advancement.
56 We also identified several lingering issues surrounding large
57 scale re-analysis.

58 Actual computation on these large datasets is a non-trivial
59 issue, as it requires access to facilities with sufficiently large,
60 high-memory machines. Amazon Web Service instances and

2 Individual components of the project are assigned specific DOIs, for ex-
ample: translated peptide files: <https://doi.org/10.5281/zenodo.745633>;
gff3 annotation files: <https://doi.org/10.5281/zenodo.744702>; annota-
tion tables: <https://doi.org/10.5281/zenodo.775129>; quantification files:
<https://doi.org/10.5281/zenodo.746294>.

3 It should be noted that uploading the assemblies to Zenodo was not an
automated process. New versions of files on Zenodo must be manually
curated. Since the start of this project, the Open Science Framework
(OSF) and the accompanying automated command-line client, osfclient
has been established. In the future, large-scale projects such as the as-
semblies created in this analysis may benefit from the integration of OSF
command-line client by automatically uploading data products to an OSF
project, which generate an OSF-specific DOI.

1 As of 17 May 2018.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

other “cloud” platforms, including XSEDE, provide flexible computing options, and are broadly accessible. Cloud-based systems, however, tend to be more expensive per computation hour than local resources. High Performance Computing (HPC) resources at local institutions represent another potential site of compute ability. However, HPCs can be temperamental and potentially balk at larger, more node-consuming procedures; moreover, bioinformatics tools may be poorly optimized for HPCs: Trinity, used in our pipeline, creates many small files for each run, and this repeatedly caused disk slowdowns on our HPC. The re-analysis by Johnson et al. [6] attempted to use both but ultimately found that the HPC provided the most consistent scalable automation for running hundreds of jobs in a cost efficient manner. However, more generally, we see no global solution for identifying and optimizing the global scientific cyberinfrastructure requirements for projects which require significant scaling; such considerations must be made on a project-by-project basis given the resources available to each lab.

Beyond the optimization of computational resources, we feel that there is a significant opportunity for scientific advancement with high-throughput sequencing projects in making data products “forward discoverable”, because this makes it possible to improve downstream work without significant upstream investment. In an ideal future, a researcher might be automatically notified when a dataset that she is actively working on is updated or changes. This presents many social and technical challenges that will need to be solved if we are to take full advantage of public datasets.

Declarations

Competing interests

The authors declare that they have no competing interests.

Funding

Funding was provided from the Gordon and Betty Moore Foundation under award number GBMF4551 to C.T.B.

Author contribution

Conceptualized by H.A., L.K.J., and C.T.B. Written by H.A. and C.T.B. Edited and revised by H.A., C.T.B., and L.K.J. All authors read and approved the final manuscript.

References

1. Caron DA, Alexander H, Allen AE, Archibald JM, Armbrust EV, Bachy C, et al. Probing the evolution, ecology and physiology of marine protists using transcriptomics. *Nature Reviews Microbiology* 2016 nov;15(1):6–20. <http://www.nature.com/doi/10.1038/nrmicro.2016.160>.
2. Keeling PJ, Burki F, Wilcox HM, Allam B, Allen EE, Amaral-Zettler LA, et al. The Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the functional diversity of eukaryotic life in the oceans through transcriptome sequencing. *PLoS Biology* 2014 jun;12(6):e1001889. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4068987>.
3. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. ABySS: A parallel assembler for short read sequence data. *Genome Research* 2009 jun;19(6):1117–1123. <http://www.ncbi.nlm.nih.gov/pubmed/19251739>.
4. Huang X, Madan A. CAP3: A DNA sequence assembly program. *Genome research* 1999 sep;9(9):868–787. <http://www.ncbi.nlm.nih.gov/pubmed/10508846>.
5. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature biotechnology* 2011 may;29(7):644–52. <http://www.ncbi.nlm.nih.gov/pubmed/21572440>.
6. Johnson LK, Alexander H, Brown CT. Re-assembly, quality evaluation, and annotation of 678 microbial eukaryotic reference transcriptomes. *GigaScience* in press;
7. MacManes MD. The Oyster River Protocol: a multi-assembler and kmer approach for de novo transcriptome assembly. *PeerJ* 2018 aug;6:e5428. <https://peerj.com/articles/5428>.
8. Shumway M, Cochrane G, Sugawara H. Archiving next generation sequencing data. *Nucleic Acids Research* 2009 jan;38:D870–D871. <http://www.ncbi.nlm.nih.gov/pubmed/19965774>.
9. Bechhofer S, Buchan I, De Roure D, Missier P, Ainsworth J, Bhagat J, et al. Why linked data is not enough for scientists. *Future Generation Computer Systems* 2013;29(2):599–611.
10. Callahan A, Cruz-Toledo J, Ansell P, Dumontier M. Bio2RDF Release 2: Improved Coverage, Interoperability and Provenance of Life Science Linked Data. In: *The Semantic Web: Semantics and Big Data* Springer, Berlin, Heidelberg; 2013.p. 200–212. http://link.springer.com/10.1007/978-3-642-38288-8_14.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65