

# Reproducible, Interactive, Scalable and Extensible Microbiome Data Science using QIIME 2

Evan Bolyen<sup>1,80</sup>, Jai Ram Rideout<sup>1,80</sup>, Matthew R. Dillon<sup>1,80</sup>, Nicholas A. Bokulich<sup>1,80</sup>, Christian C. Abnet<sup>2</sup>, Gabriel A. Al-Ghalith<sup>3</sup>, Harriet Alexander<sup>4,5</sup>, Eric J. Alm<sup>6,7</sup>, Manimozhayan Arumugam<sup>8</sup>, Francesco Asnicar<sup>9</sup>, Yang Bai<sup>10,11,12</sup>, Jordan E. Bisanz<sup>13</sup>, Kyle Bittinger<sup>14,15</sup>, Asker Brejnrod<sup>8</sup>, Colin J. Brislawn<sup>16</sup>, C. Titus Brown<sup>5</sup>, Benjamin J. Callahan<sup>17,18</sup>, Andrés Mauricio Caraballo-Rodríguez<sup>19</sup>, John Chase<sup>1</sup>, Emily K. Cope<sup>1,20</sup>, Ricardo Da Silva<sup>19</sup>, Christian Diener<sup>21</sup>, Pieter C. Dorrestein<sup>19</sup>, Gavin M. Douglas<sup>22</sup>, Daniel M. Durall<sup>23</sup>, Claire Duvallet<sup>6</sup>, Christian F. Edwardson<sup>24</sup>, Madeleine Ernst<sup>19,25</sup>, Mehrbod Estaki<sup>26</sup>, Jennifer Fouquier<sup>27,28</sup>, Julia M. Gauglitz<sup>19</sup>, Sean M. Gibbons<sup>21,29</sup>, Deanna L. Gibson<sup>30,31</sup>, Antonio Gonzalez<sup>32</sup>, Kestrel Gorlick<sup>1</sup>, Jiarong Guo<sup>33</sup>, Benjamin Hillmann<sup>34</sup>, Susan Holmes<sup>35</sup>, Hannes Holste<sup>32,36</sup>, Curtis Huttenhower<sup>37,38</sup>, Gavin A. Huttley<sup>39</sup>, Stefan Janssen<sup>40</sup>, Alan K. Jarmusch<sup>19</sup>, Lingjing Jiang<sup>41</sup>, Benjamin D. Kaehler<sup>39,42</sup>, Kyo Bin Kang<sup>19,43</sup>, Christopher R. Keefe<sup>1</sup>, Paul Keim<sup>1</sup>, Scott T. Kelley<sup>44</sup>, Dan Knights<sup>34,45</sup>, Irina Koester<sup>19,46</sup>, Tomasz Kosciolk<sup>47</sup>, Jorden Kreps<sup>1</sup>, Morgan G. I. Langille<sup>48</sup>, JoslynnLee<sup>49</sup>, Ruth Ley<sup>50,51</sup>, Yong-Xin Liu<sup>10,11</sup>, Erika Lofffield<sup>2</sup>, Catherine Lozupone<sup>28</sup>, Massoud Maher<sup>52</sup>, Clarisse Marotz<sup>32</sup>, Bryan D. Martin<sup>53</sup>, Daniel McDonald<sup>32</sup>, Lauren J. McIver<sup>37,38</sup>, Alexey V. Melnik<sup>19</sup>, Jessica L. Metcalf<sup>54</sup>, Sydney C. Morgan<sup>55</sup>, Jamie T. Morton<sup>32,52</sup>, Ahmad Turan Naimey<sup>1</sup>, Jose A. Navas- Molina<sup>32,52,56</sup>, Louis Felix Nothias<sup>19</sup>, Stephanie B. Orchanian<sup>57</sup>, Talima Pearson<sup>1</sup>, Samuel L. Peoples<sup>58,59</sup>, Daniel Petras<sup>19</sup>, Mary Lai Preuss<sup>60</sup>, Elmar Pruesse<sup>28</sup>, Lasse Buur Rasmussen<sup>8</sup>, Adam Rivers<sup>61</sup>, Michael S. Robeson II<sup>62</sup>, Patrick Rosenthal<sup>60</sup>, Nicola Segata<sup>9</sup>, Michael Shaffer<sup>27,28</sup>, Arron Shiffer<sup>1</sup>, Rashmi Sinha<sup>2</sup>, Se Jin Song<sup>32</sup>, John R. Spear<sup>63</sup>, Austin D. Swafford<sup>57</sup>, Luke R. Thompson<sup>64,65</sup>, Pedro J. Torres<sup>66</sup>, Pauline Trinh<sup>67</sup>, Anupriya Tripathi<sup>19,32,68</sup>, Peter J. Turnbaugh<sup>69</sup>, Sabah Ul-Hasan<sup>70</sup>, Justin J. J. van der Hooft<sup>71</sup>, Fernando Vargas<sup>68</sup>, Yoshiki Vázquez-Baeza<sup>32</sup>, Emily Vogtmann<sup>2</sup>, Max von Hippel<sup>72</sup>, William Walters<sup>50</sup>, Yunhu Wan<sup>2</sup>, Mingxun Wang<sup>19</sup>, Jonathan Warren<sup>73</sup>, Kyle C. Weber<sup>61,74</sup>, Charles H. D. Williamson<sup>75</sup>, Amy D. Willis<sup>76</sup>, Zhenjiang Zech Xu<sup>32</sup>, Jesse R. Zaneveld<sup>77</sup>, Yilong Zhang<sup>78</sup>, Qiyun Zhu<sup>32</sup>, Rob Knight<sup>32,57,79</sup>, J. Gregory Caporaso<sup>1,20,\*</sup>

For publication in *Nature Biotechnology*

## Author Affiliations

- <sup>1</sup>Center for Applied Microbiome Science, Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA
- <sup>2</sup>Metabolic Epidemiology Branch, National Cancer Institute, Rockville, MD, USA
- <sup>3</sup>Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA
- <sup>4</sup>Biology Department, Woods Hole Oceanographic Institution, Woods Hole, MA, USA
- <sup>5</sup>Department of Population Health and Reproduction, University of California, Davis, Davis, CA, USA
- <sup>6</sup>Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA
- <sup>7</sup>Center for Microbiome Informatics and Therapeutics, Massachusetts Institute of Technology, Cambridge, MA, USA
- <sup>8</sup>Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark
- <sup>9</sup>Centre for Integrative Biology, University of Trento, Trento, Italy
- <sup>10</sup>State Key Laboratory of Plant Genomics, Institute of Genetics and Developmental Biology, Chinese Academy of Sciences, Beijing, China
- <sup>11</sup>Centre of Excellence for Plant and Microbial Sciences (CEPAMS), Institute of Genetics and Developmental Biology, Chinese Academy of Sciences & John Innes Centre, Beijing, China
- <sup>12</sup>University of Chinese Academy of Sciences, Beijing, China
- <sup>13</sup>Department of Microbiology and Immunology, University of California, San Francisco, San Francisco, CA, USA
- <sup>14</sup>Division of Gastroenterology and Nutrition, Children's Hospital of Philadelphia, Philadelphia, PA, USA
- <sup>15</sup>Hepatology, Children's Hospital of Philadelphia, Philadelphia, PA, USA
- <sup>16</sup>Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, WA, USA
- <sup>17</sup>Department of Population Health & Pathobiology, North Carolina State University, Raleigh, NC, USA
- <sup>18</sup>Bioinformatics Research Center, North Carolina State University, Raleigh, NC, USA
- <sup>19</sup>Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, San Diego, CA, USA
- <sup>20</sup>Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA
- <sup>21</sup>Institute for Systems Biology, Seattle, WA, USA
- <sup>22</sup>Department of Microbiology and Immunology, Dalhousie University, Halifax, Nova Scotia, Canada
- <sup>23</sup>Irving K. Barber School of Arts and Sciences, University of British Columbia, Kelowna, British Columbia, Canada
- <sup>24</sup>A. Watson Armour III Center for Animal Health and Welfare, Aquarium Microbiome Project, John G. Shedd Aquarium, Chicago, IL, USA
- <sup>25</sup>Department of Congenital Disorders, Statens Serum Institut, Copenhagen, Denmark
- <sup>26</sup>Department of Biology, University of British Columbia Okanagan, Okanagan, British Columbia, Canada

- <sup>27</sup>Computational Bioscience Program, University of Colorado Anschutz Medical Campus, Aurora, CO, USA
- <sup>28</sup>Department of Medicine, Division of Biomedical Informatics and Personalized Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO, USA
- <sup>29</sup>eScience Institute, University of Washington, Seattle, WA, USA
- <sup>30</sup>Irving K. Barber School of Arts and Sciences, Department of Biology, University of British Columbia, Kelowna, British Columbia, Canada
- <sup>31</sup>Department of Medicine, University of British Columbia, Kelowna, British Columbia, Canada
- <sup>32</sup>Department of Pediatrics, University of California San Diego, La Jolla, CA, USA
- <sup>33</sup>Center for Microbial Ecology, Michigan State University, East Lansing, MI, USA
- <sup>34</sup>Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA
- <sup>35</sup>Statistics Department, Stanford University, Palo Alto, CA, USA
- <sup>36</sup>Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA
- <sup>37</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA
- <sup>38</sup>Broad Institute of MIT and Harvard, Cambridge, MA, USA
- <sup>39</sup>Research School of Biology, The Australian National University, Canberra, Australian Capital Territory, Australia
- <sup>40</sup>Department of Pediatric Oncology, Hematology and Clinical Immunology, Heinrich-Heine University Dusseldorf, Dusseldorf, Germany
- <sup>41</sup>Department of Family Medicine and Public Health, University of California San Diego, La Jolla, CA, USA
- <sup>42</sup>School of Science, University of New South Wales, Canberra, Australian Capital Territory, Australia
- <sup>43</sup>College of Pharmacy, Sookmyung Women's University, Seoul, Republic of Korea
- <sup>44</sup>Department of Biology, San Diego State University, San Diego, CA, USA
- <sup>45</sup>Biotechnology Institute, University of Minnesota, Saint Paul, MN, USA
- <sup>46</sup>Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA
- <sup>47</sup>Department of Pediatrics, University of California San Diego, La Jolla, California, USA
- <sup>48</sup>Department of Pharmacology, Dalhousie University, Halifax, Nova Scotia, Canada
- <sup>49</sup>Science Education, Howard Hughes Medical Institute, Ashburn, VA, USA
- <sup>50</sup>Department of Microbiome Science, Max Planck Institute for Developmental Biology, Tübingen, Germany
- <sup>51</sup>Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY, USA
- <sup>52</sup>Department of Computer Science & Engineering, University of California San Diego, La Jolla, CA, USA
- <sup>53</sup>Department of Statistics, University of Washington, Seattle, WA, USA
- <sup>54</sup>Department of Animal Science, Colorado State University, Fort Collins, CO, USA

- <sup>55</sup>Irving K. Barber School of Arts and Sciences, Unit 2 (Biology), University of British Columbia, Kelowna, British Columbia, Canada
- <sup>56</sup>Google LLC, Mountain View, CA, USA
- <sup>57</sup>Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, USA
- <sup>58</sup>School of Information Studies, Syracuse University, Syracuse, NY, USA
- <sup>59</sup>School of STEM, University of Washington Bothell, Bothell, WA, USA
- <sup>60</sup>Department of Biological Sciences, Webster University, St. Louis, MO, USA
- <sup>61</sup>Agricultural Research Service, Genomics and Bioinformatics Research Unit, United States Department of Agriculture, Gainesville, FL, USA
- <sup>62</sup>College of Medicine, Department of Biomedical Informatics, University of Arkansas for Medical Sciences, Little Rock, AR, USA
- <sup>63</sup>Department of Civil and Environmental Engineering, Colorado School of Mines, Golden, CO, USA
- <sup>64</sup>Department of Biological Sciences and Northern Gulf Institute, University of Southern Mississippi, Hattiesburg, MS, USA
- <sup>65</sup>Ocean Chemistry and Ecosystems Division, Atlantic Oceanographic and Meteorological Laboratory, National Oceanic and Atmospheric Administration, Miami, FL, USA
- <sup>66</sup>Department of Biology, San Diego State University, San Diego, CA, USA
- <sup>67</sup>Department of Environmental and Occupational Health Sciences, University of Washington, Seattle, WA, USA
- <sup>68</sup>Division of Biological Sciences, University of California San Diego, San Diego, CA, USA
- <sup>69</sup>Department of Microbiology and Immunology, University of California San Francisco, San Francisco, CA, USA
- <sup>70</sup>Quantitative and Systems Biology Graduate Program, University of California Merced, Merced, CA, USA
- <sup>71</sup>Bioinformatics Group, Wageningen University, Wageningen, the Netherlands
- <sup>72</sup>Department of Mathematics, University of Arizona, Tucson, AZ, USA
- <sup>73</sup>National Laboratory Service, Environment Agency, Starcross, UK
- <sup>74</sup>College of Agriculture and Life Sciences, University of Florida, Gainesville, FL, USA
- <sup>75</sup>Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA
- <sup>76</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA
- <sup>77</sup>School of STEM, Division of Biological Sciences, University of Washington Bothell, Bothell, WA, USA
- <sup>78</sup>Merck & Co. Inc., Kenilworth, NJ, USA
- <sup>79</sup>Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA
- <sup>80</sup>These authors contributed equally: Evan Bolyen, Jai Ram Rideout, Matthew R. Dillon, Nicholas A. Bokulich

**To the Editor** — Rapid advances in DNA-sequencing and bioinformatics technologies in the past two decades have substantially improved understanding of the microbial world. This growing understanding relates to the vast diversity of microorganisms; how microbiota and microbiomes affect disease<sup>1</sup> and medical treatment<sup>2</sup>; how microorganisms affect the health of the planet<sup>3</sup>; and the nascent exploration of the medical<sup>4</sup>, forensic<sup>5</sup>, environmental<sup>6</sup> and agricultural<sup>7</sup> applications of microbiome biotechnology. Much of this work has been driven by marker-gene surveys (for example, bacterial/archaeal 16S rRNA genes, fungal internal-transcribed-spacer regions and eukaryotic 18S rRNA genes), which profile microbiota with varying degrees of taxonomic specificity and phylogenetic information. The field is now transitioning to integrate other data types, such as metabolite<sup>8</sup>, metaproteome<sup>9</sup> or metatranscriptome<sup>9,10</sup> profiles.

The QIIME 1 microbiome bioinformatics platform has supported many microbiome studies and gained a broad user and developer community. Interactions with QIIME 1 users in our online support forum, our workshops and direct collaborations have shown the platform's potential to serve an increasingly diverse array of microbiome researchers in academia, government and industry. Here, we present QIIME 2, a completely reengineered and rewritten system that is expected to facilitate reproducible and modular analysis of microbiome data to enable the next generation of microbiome science.

QIIME 2 was developed on the basis of a plugin architecture (Supplementary Fig. 1) that allows third parties to contribute functionality (<https://library.qiime2.org>). QIIME 2 plugins exist for latest-generation tools for sequence quality control from different sequencing platforms (DADA2 (ref. 11) and Deblur<sup>12</sup>), taxonomy assignment<sup>13</sup> and phylogenetic insertion<sup>14</sup>, which quantitatively improve the results over QIIME 1 and other tools (as detailed in the corresponding tool-specific publications). The plugins also support qualitatively new functionality, including microbiome paired-sample and time-series analysis<sup>15</sup> (which are critical for studying the effects of treatments on the microbiome), and machine learning<sup>16</sup>. Trained machine learning models can be saved for application to new data and interrogated to identify important microbiome features. Several recently released plugins, including q2-cscs<sup>17</sup>, q2-metabolomics<sup>18</sup>, q2-shogun<sup>19</sup>, q2-metaphlan2 (ref. 20) and q2-picrust2 (ref. 21), provide initial support for analysis of metabolomics and shotgun metagenomics data. We are currently working with teams developing bioinformatics tools for metatranscriptomics and metaproteomics, and we expect to

add new plugins supporting these data types to the ecosystem shortly. Additionally, many of the existing 'downstream' analysis tools, such as q2-sample-classifier<sup>16</sup>, can already work with these data types individually or in combination if they are provided in a feature table. Thus, QIIME 2 has the potential to serve not only as a marker-gene analysis tool but also a multidimensional and powerful data science platform that can be rapidly adapted to analyze diverse microbiome features.

QIIME 2 provides many new interactive visualization tools facilitating exploratory analyses and result reporting. Static versions of interactive visualizations resulting from four worked examples are provided in Fig. 1. QIIME 2 View (<https://view.qiime2.org>) is a unique new service (Supplementary Methods) that allows users to securely share and interact with results without installing QIIME 2. The QIIME 2 visualizations presented in Fig. 1 are provided in Supplementary File 1 to allow readers to interact with QIIME 2 View. Corresponding worked QIIME 2 example code is provided in the Supplementary Methods.

Reproducibility, transparency and clarity of microbiome data science are guiding principles in QIIME 2 design. To this end, QIIME 2 includes a decentralized data-provenance tracking system: details of all analysis steps with references to intermediate data are automatically stored in the results. Users can thus retrospectively determine exactly how any result was generated (Fig. 2 illustrates a simplified provenance graph derived from the data provenance of Fig. 1b). QIIME 2 also detects corrupted results indicating that the provenance is no longer reliable and the results no longer contain information enabling reproducibility. The provenance of the visualizations presented in Fig. 1 can be interactively reviewed by loading the contents of Supplementary File 1 with QIIME 2 View, providing far more detailed information than can typically be provided in Methods text. QIIME 2 results are also semantically typed (Fig. 2), and actions indicate acceptable input types, clarifying the data that actions should be applied to and making complex workflows less error prone. Complex workflows can be created and shared by using Jupyter Notebooks<sup>22</sup> or Common Workflow Language (CWL)<sup>23</sup>, and support for other workflow engines is currently in development.

Finally, QIIME 2 provides a software-development kit (<https://dev.qiime2.org>) that can be used to integrate it as a component of other systems (such as Qiita<sup>24</sup> or Illumina BaseSpace) and to develop interfaces targeted toward users with different levels of computational sophistication (Supplementary Fig. 2). QIIME 2 provides the QIIME 2 Studio graphical user interface and QIIME 2 View, interfaces designed for end-user biologists, clinicians and policy-

makers; the QIIME 2 application programming interface, designed for data scientists who want to automate workflows or work interactively in Jupyter Notebooks<sup>22</sup>; and q2cli and q2cwl, providing a command-line interface and CWL<sup>23</sup> wrappers for QIIME 2, designed for experts in high-performance computing. At present, computationally expensive steps support parallel computing at the individual-action level (for example, many actions including de-noising and taxonomy assignment support multiple threads). We are currently developing deeper integration with parallelism strategies available in third-party workflow engines, and workflow-level parallelism is currently possible through CWL.

There are many other powerful open-source software tools for microbiome data science, including mothur<sup>25</sup>, phyloseq<sup>26</sup> and related tools available through Bioconductor<sup>27</sup>, and the biobakery suite<sup>20,21,28</sup>. The microbiome bioinformatics platform mothur is often compared to QIIME 1 and QIIME 2. A major difference between mothur and QIIME lies in the interactive visualizations: QIIME 2 provides many interactive visualization tools (several examples are provided in Fig. 1), whereas mothur focuses on generating data that can be easily loaded and visualized with other tools. The phyloseq tool focuses on microbiome statistical analysis and generating publication-ready visualizations but, unlike QIIME 2, begins with a feature or operational-taxonomic-unit table, leaving ‘upstream’ processing steps, such as sequence demultiplexing and quality control, to other processing pipelines, many of which (like phyloseq) are available through Bioconductor. The biobakery suite provides analytic functionality that complements that of QIIME 2, and we are actively working with biobakery developers to support interoperability by making their tools accessible as QIIME 2 plugins (for example, the q2-metaphlan2 plugin allows users to run MetaPhlan2 through QIIME 2). QIIME 2 provides the only Python-based microbiome data-science platform that supports retrospective data-provenance tracking to ensure reproducibility, multi-omics analysis support, interfaces geared toward different user types to enhance usability and an extensibility-focused design through the plugin architecture and software-development kit. We share feedback from users of QIIME 2 on these and other features in Supplementary Methods.

The tools described in the preceding paragraph are all interoperable through plugins, exchange of files in standard formats or using multi-language environments, such as Jupyter Notebooks<sup>22</sup>. For example, the BIOM format<sup>29</sup> is supported by all of them. A diverse ecosystem of interoperable software is beneficial for the field, because it allows both experienced users to obtain multiple perspectives on their data and novice bioinformaticians to work in the

programming environments that they are most comfortable with (for example, phyloseq allows users to work in R, whereas QIIME 2 allows users to work in Python). We plan to continue working with the developers of these tools, and with organizations such as the Genomics Standards Consortium, on plugins and standards to ensure interoperability, as well as developing tools to automatically import data from microbiome data-sharing platforms such as Qiita, the European Bioinformatics Institute (EBI) European Read Archive and the National Center for Biotechnology Information (NCBI) Sequence Read Archive.

Advances in microbiome research promise to improve many aspects of health and the world, and QIIME 2 will help drive those advances by enabling accessible, community-driven microbiome data science.

### **Data availability**

Data for the analyses presented in Fig. 1 are available as follows: Earth Microbiome Project data in Fig. 1a were obtained from <ftp://ftp.microbio.me/emp/release1>, and the American Gut Project (AGP) data were obtained from Qiita (<http://qiita.microbio.me>) study ID 10317. Sequence data in Fig. 1c are available in Qiita under study ID 10249 and the EBI under accession number ERP016173. Sequence data in Fig. 1b are available in Qiita under study ID 925 and the EBI under accession number ERP022167. Data in Fig. 1d are available in the q2-ili GitHub repository (<https://github.com/biocore/q2-ili>). Interactive versions of the Fig. 1 visualizations can be accessed at <https://github.com/qiime2/paper1>.

### **Code availability**

QIIME 2 is open source and free for all use, including commercial. It is licensed under a BSD three-clause license. Source code is available at <https://github.com/qiime2>. Help for QIIME 2 is provided at <https://forum.qiime2.org>.

### **Supplementary Material**

Refer to Web version on PubMed Central for supplementary material.

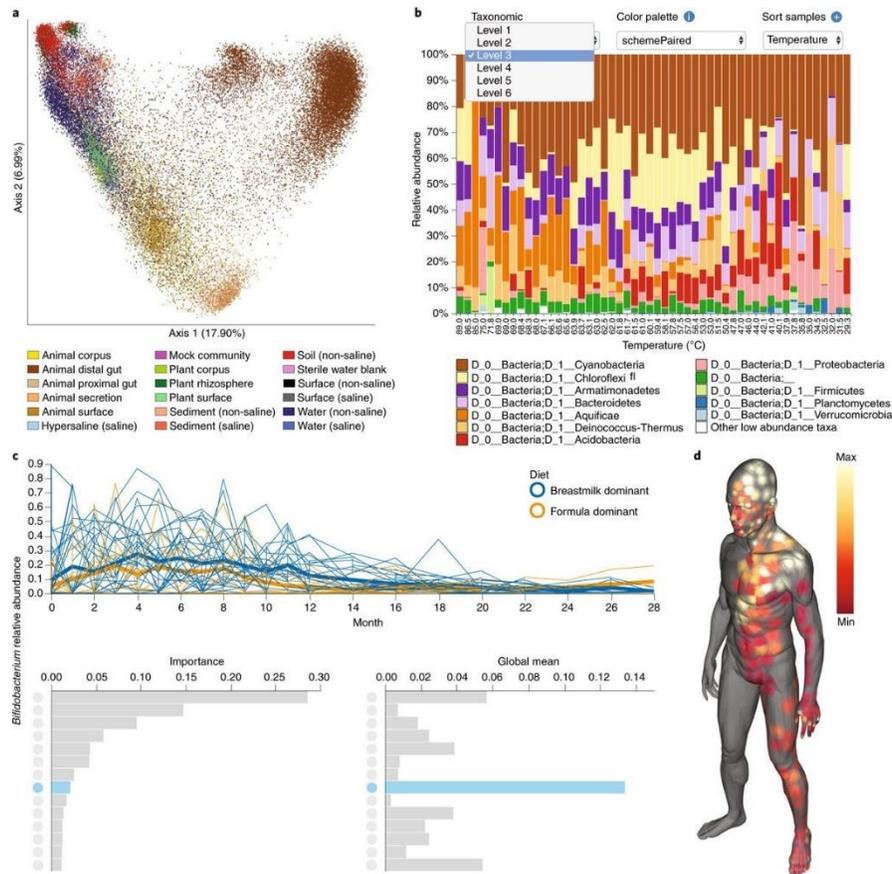
## Acknowledgements

QIIME 2 development was primarily funded by NSF Awards 1565100 to J.G.C. and 1565057 to R.K. Partial support was also provided by the following: grants NIH U54CA143925 (J.G.C. and T.P.) and U54MD012388 (J.G.C. and T.P.); grants from the Alfred P. Sloan Foundation (J.G.C. and R.K.); ERCSTG project MetaPG (N.S.); the Strategic Priority Research Program of the Chinese Academy of Sciences QYZDB-SSW-SMC021 (Y.B.); the Australian National Health and Medical Research Council APP1085372 (G.A.H., J.G.C., Von Bing Yap and R.K.); the Natural Sciences and Engineering Research Council (NSERC) to D.L.G.; and the State of Arizona Technology and Research Initiative Fund (TRIF), administered by the Arizona Board of Regents, through Northern Arizona University. All NCI coauthors were supported by the Intramural Research Program of the National Cancer Institute.

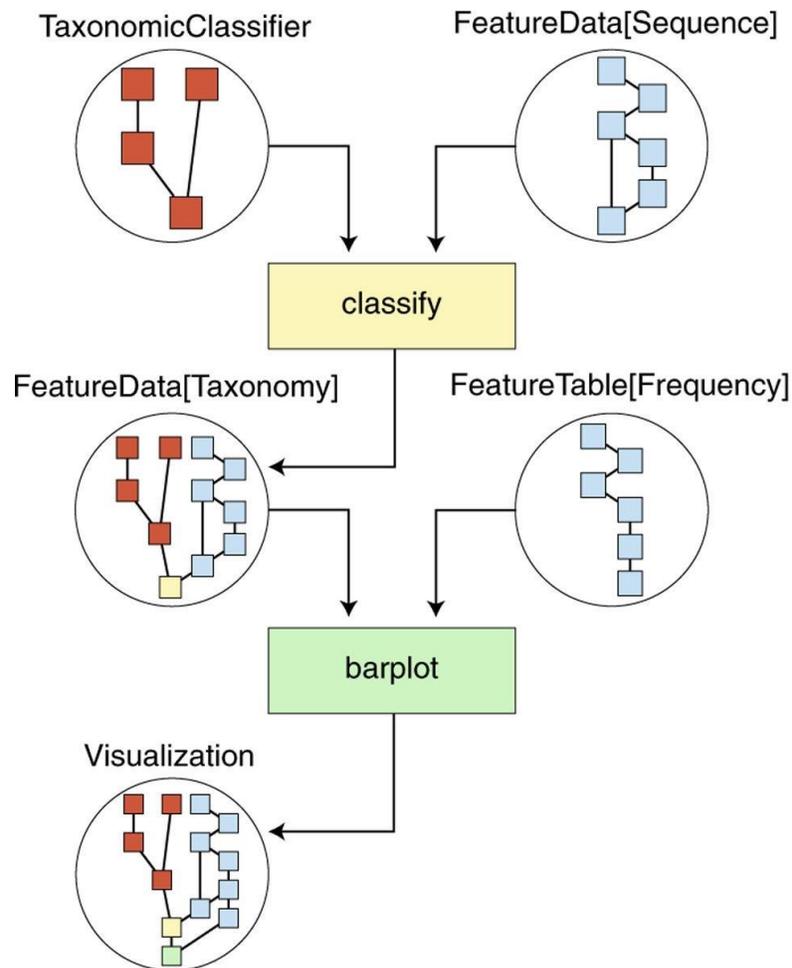
S.M.G. and C. Diener were supported by the Washington Research Foundation Distinguished Investigator Award. Thanks to the Yellowstone Center for Resources for research permit no. 5664 to J.R.S. for Yellowstone access and sample collection. We thank P. J. McMurdie for helpful discussion on the relationships between QIIME 2 and phyloseq. We would like to thank the users of QIIME 1 and 2, whose invaluable feedback has shaped QIIME 2. In particular, we would like to thank A. Abdelfattah (Stockholm University, Sweden), R. C. T. Boutin (University of British Columbia, Canada), D. J. Bradshaw II (Florida Atlantic University Harbor Branch Oceanographic Institute, USA), L. Bullington (MPG Ranch, USA), J. W. Debelius (Karolinska Institutet, Sweden), C. Duvallet (Massachusetts Institute of Technology, USA), E. Korzune Ganda (Cornell University, USA), A. Mahnert (Medical University of Graz, Austria), M. C. Melendrez (St. Cloud State University, USA), D. O'Rourke (University of New Hampshire, USA), A. R. Rivers (USDA ARS, USA), B. Sen (Tianjin University, China), S. Tangedal (Haukeland University Hospital and University of Bergen, Norway), P. J. Torres (San Diego State University, USA) and J. Warren (National Laboratory Service, UK) for writing end-user reviews included in the Supplementary Methods.

## References

1. Smith MI et al. *Science* 339, 548–554 (2013).
2. Gopalakrishnan V et al. *Science* 359, 97–103 (2018).
3. Gehring CA, Sthultz CM, Flores-Rentería L, Whipple AV & Whitham TG *Proc. Natl Acad. Sci. USA* 114, 11169–11174 (2017).
4. Lee K, Pletcher SD, Lynch SV, Goldberg AN & Cope EK *Front. Cell. Infect. Microbiol.* 8, 168(2018).
5. Metcalf JL et al. *Science* 351, 158–162 (2016).
6. Rubin RL et al. *Ecol. Appl.* 28, 1594–1605 (2018).
7. Pineda A, Kaplan I & Bezemer TM *Trends Plant Sci.* 22, 770–778 (2017).
8. Kapono CA et al. *Sci. Rep.* 8, 3669 (2018).
9. Verberkmoes NC et al. *ISME J.* 3, 179–189 (2009).
10. Barr T et al. *Gut Microbes* 9, 338–356 (2018).
11. Callahan BJ et al. *Nat. Methods* 13, 581–3 (2016).
12. Amir A et al. *mSystems* 2, e00191–16 (2017).
13. Bokulich NA et al. *Microbiome* 6, 90 (2018).
14. Janssen S et al. *mSystems* 3, e00021–18 (2018).
15. Bokulich NA et al. *mSystems* 3, e00219–18 (2018).
16. Bokulich N et al. *J. Open Source Softw.* 3, 934 (2018).
17. Sedio BE, Rojas Echeverri JC, Boya PCA & Wright SJ *Ecology* 98, 616–623 (2017).
18. Wang M et al. *Nat. Biotechnol.* 34, 828–837 (2016).
19. Hillmann B et al. *mSystems.* 3, e00069–18 (2018).
20. Truong DT et al. *Nat. Methods* 12, 902–903 (2015).
21. Langille MGI et al. *Nat. Biotechnol.* 31, 814–821 (2013).
22. Kluyver T et al. Positioning and power in academic publishing: players, agents and agendas. in *Proc 20th International Conference on Electronic Publishing* (eds Loizides F & Schmidt B) 87–90 (IOS Press, 2016).
23. Amstutz P et al. 10.6084/m9.figshare.3115156.v2 (2016).
24. Gonzalez A et al. *Nat. Methods* 15, 796–798 (2018).
25. Schloss PD et al. *Appl. Environ. Microbiol.* 75, 7537–7541 (2009).
26. McMurdie PJ & Holmes S *PLoS One* 8, e61217 (2013).
27. 27. Huber W et al. *Nat. Methods* 12, 115–121 (2015).
28. Franzosa EA et al. *Nat. Methods* 15, 962–968 (2018).
29. McDonald D et al. *Gigascience* 1, 7 (2012).



**Figure 1. QIIME 2 provides many interactive visualization tools.** The products of four worked examples are presented here, and interactive versions of these screen captures are available in Supplementary File 1 and at <https://github.com/qiime2/paper1>. Detailed descriptions and methods, including the commands used to generate each of these visualizations, are provided in Supplementary Methods. **a**, Unweighted UniFrac principal coordinate analysis plot containing 37,680 samples, illustrating the scalability of QIIME 2. Colors indicate sample type, as described by the Earth Microbiome Project ontology (EMPO). **b**, Interactive taxonomic composition bar plot illustrating the phylum- level composition of microbial-mat samples collected along a temperature gradient in Yellowstone National Park Hot Spring outflow channels (Steep Cone Geyser). The many interactive controls available in this plot vastly decrease the burden of exploratory analysis over QIIME 1. **c**, Feature volatility plot (<https://msystems.asm.org/content/3/6/e00219-18>) illustrating the change in *Bifidobacterium* abundance over time in breast-fed and formula- fed infants. Temporally interesting features can be interactively discovered with this visualization. Bar charts rank the importance (predictive power for time point) and mean abundance of all microbial features. These bar charts provide an interface for visualizing volatility plots (line plots) of individual features in the context of their importance and abundance; clicking on a bar will display the volatility plot of that feature and highlight in blue that feature's importance and abundance in the bar charts below. **d**, Molecular cartography of the human skin surface. Colored spots represent the abundance of the small- molecule cosmetic ingredient sodium laureth sulfate on the human skin. Sample data can be interactively visualized in three-dimensional models, thus supporting the discovery of spatial patterns.



**Figure 2. QIIME 2 iteratively records data provenance, ensuring bioinformatics reproducibility.**

This simplified diagram illustrates the automatically tracked information regarding the creation of the taxonomy bar plot presented in Fig. 1b. QIIME 2 results (circles) contain network diagrams illustrating the data provenance stored in the result. Actions (quadrilaterals) are applied to QIIME 2 results and generate new results. Arrows indicate the flow of QIIME 2 results through actions.

TaxonomicClassifier and FeatureData[Sequence] inputs contain independent provenance (red and blue, respectively) and are provided to a classify action (yellow), which taxonomically annotates sequences. The result of the classify action, a FeatureData[Taxonomy] result, integrates the provenance of both inputs with the classify action. This result is then provided to the barplot action with a FeatureTable[Frequency] input, which shares some provenance with the FeatureData[Sequence] input, because they were generated from the same upstream analysis. The resulting visualization (Fig. 1b) has the complete data provenance and correctly identifies shared processing of inputs. This simplified representation was created manually from the complete provenance graph for the purpose of illustration. An interactive and complete version of this provenance graph (as well as those for other Fig. 1 panels) can be accessed through Supplementary File 1.