

2 QIIME 2: Reproducible, interactive, scalable, and extensible microbiome data science

4 Evan Bolyen^{1,*}, Jai Ram Rideout^{1,*}, Matthew R Dillon^{1,*}, Nicholas A Bokulich^{1,*}, Christian C Abnet², Gabriel A
Al-Ghalith³, Harriet Alexander^{4,5}, Eric J Alm^{6,7}, Manimozhiyan Arumugam⁸, Francesco Asnicar⁹, Yang Bai^{10,11,12},
6 Jordan E Bisanz¹³, Kyle Bittinger^{14,15}, Asker Brejnrod¹⁶, Colin J Brislawn¹⁷, C Titus Brown⁵, Benjamin J
Callahan^{18,19}, Andrés Mauricio Caraballo-Rodríguez²⁰, John Chase¹, Emily K Cope^{1,21}, Ricardo Da Silva²⁰,
8 Pieter C Dorrestein²⁰, Gavin M Douglas²², Daniel M Durall²³, Claire Duvallet⁶, Christian F Edwardson²⁴,
Madeleine Ernst²⁰, Mehrbod Estaki²⁵, Jennifer Fouquier^{26,27}, Julia M Gauglitz²⁰, Deanna L Gibson^{28,29}, Antonio
Gonzalez³⁰, Kestrel Gorlick¹, Jiarong Guo³¹, Benjamin Hillmann³², Susan Holmes³³, Hannes Holste^{30,34}, Curtis
10 Huttenhower^{35,36}, Gavin A Huttley³⁷, Stefan Janssen³⁸, Alan K Jarmusch²⁰, Lingjing Jiang³⁹, Benjamin D
Kaehler³⁷, Kyo Bin Kang^{40,20}, Christopher R Keefe¹, Paul Keim¹, Scott T Kelley⁴¹, Dan Knights^{42,32}, Irina
12 Koester^{43,20}, Tomasz Kosciolk⁴⁴, Jordan Kreps¹, Morgan GI Langille⁴⁵, Joslynn Lee⁴⁶, Ruth Ley^{47,48}, Yong-Xin
Liu^{10,11}, Erika Lofffield², Catherine Lozupone⁴⁹, Massoud Maher⁵⁰, Clarisse Marotz³⁰, Bryan D Martin⁵¹, Daniel
14 McDonald³⁰, Lauren J McIver^{35,36}, Alexey V Melnik²⁰, Jessica L Metcalf⁵², Sydney C Morgan⁵³, Jamie T
Morton^{30,50}, Ahmad Turan Naimey¹, Jose A Navas-Molina^{50,30,54}, Louis Felix Nothias²⁰, Stephanie B Orchanian⁵⁵,
16 Talima Pearson¹, Samuel L Peoples^{56,57}, Daniel Petras²⁰, Mary Lai Preuss⁵⁸, Elmar Pruesse⁴⁹, Lasse Buur
Rasmussen¹⁶, Adam Rivers⁵⁹, Michael S Robeson, II⁶⁰, Patrick Rosenthal⁵⁸, Nicola Segata⁹, Michael
18 Shaffer^{49,61}, Arron Shiffer¹, Rashmi Sinha², Se Jin Song³⁰, John R Spear⁶², Austin D Swafford⁵⁵, Luke R
Thompson^{63,64}, Pedro J Torres⁶⁵, Pauline Trinh⁶⁶, Anupriya Tripathi^{20,30,67}, Peter J Turnbaugh⁶⁸, Sabah
20 Ul-Hasan⁶⁹, Justin JJ van der Hoof⁷⁰, Fernando Vargas⁶⁷, Yoshiki Vázquez-Baeza³⁰, Emily Vogtmann², Max
von Hippel⁷¹, William Walters⁴⁷, Yunhu Wan², Mingxun Wang²⁰, Jonathan Warren⁷², Kyle C Weber^{59,73}, Charles
22 HD Williamson¹, Amy D Willis⁷⁴, Zhenjiang Zech Xu³⁰, Jesse R Zaneveld⁷⁵, Yilong Zhang⁷⁶, Qiyun Zhu³⁰, Rob
Knight^{30,77,55}, and J Gregory Caporaso^{1,21,+}

24 * These authors contributed equally to this work.

+ Please address correspondence to gregcaporaso@gmail.com.

26 Author affiliations are provided following the Main Text References section.

To get help with QIIME 2, visit <https://forum.qiime2.org>.

28 Abstract

30 We present QIIME 2, an open-source microbiome data science platform accessible to users spanning the
microbiome research ecosystem, from scientists and engineers to clinicians and policy makers. QIIME 2
32 provides new features that will drive the next generation of microbiome research. These include interactive
spatial and temporal analysis and visualization tools, support for metabolomics and shotgun metagenomics
analysis, and automated data provenance tracking to ensure reproducible, transparent microbiome data
34 science.

Main text

36 Rapid advances in DNA sequencing and bioinformatics technologies in the past two decades have significantly
37 improved our understanding of the microbial world. These include our growing understanding of the vast
38 diversity of microorganisms; how our microbiota and microbiomes impact disease¹ and medical treatment²;
39 how microorganisms impact the health of our planet³; and our nascent exploration of the medical⁴, forensic⁵,
40 environmental⁶, and agricultural⁷ applications of microbiome biotechnology. Much of this work has been driven
41 by marker gene surveys (e.g., bacterial/archaeal 16S rRNA genes, fungal ITS, eukaryal 18S rRNA genes),
42 which profile microbiota with varying degrees of taxonomic specificity and phylogenetic information. The field is
43 now transitioning to integrate other data types, such as metabolite⁸ or metatranscriptome⁹ profiles.

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

50 QIIME 2 is developed based on a plugin architecture (Figure S1) that allows third-parties to contribute
51 functionality (see <https://library.qiime2.org>). QIIME 2 plugins exist for latest-generation tools for sequence
52 quality control from different sequencing platforms (DADA2¹⁰ and Deblur¹¹), taxonomy assignment¹², and
53 phylogenetic insertion¹³, that quantitatively improve results over QIIME 1 and other tools (detailed in the
54 corresponding tool-specific publications). Plugins also support qualitatively new functionality including
55 microbiome paired-sample and time-series analysis¹⁴, critical for studying the impact of treatment on the
56 microbiome, and for machine learning¹⁵, including the ability to save trained models and apply them to new
57 data and to interrogate models to identify important microbiome features. Several recently released plugins,
58 including q2-cscs¹⁶, q2-metabolomics¹⁷, q2-shogun¹⁸, q2-metaphlan2¹⁹, and q2-picrust2²⁰, provide initial
59 support for analysis of metabolomics and shotgun metagenomics data. This marks the potential of QIIME 2 to
60 serve not only as a marker gene analysis tool, but also a multi-dimensional and powerful data science platform
61 that can be rapidly adapted to analyze diverse microbiome features.

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

68 Reproducibility, transparency, and clarity of microbiome data science are guiding principles in the QIIME 2
69 design. Toward this end, it includes a decentralized data provenance tracking system: details of all analysis
70 steps with references to intermediate data are automatically stored in the results. Users can thus
71 retrospectively determine exactly how any result was generated (Figure 2). QIIME 2 also detects corrupted
72 results, indicating that provenance is no longer reliable and the results no longer contain information enabling
73 reproducibility. Provenance of the visualizations presented in Figure 1 can be interactively reviewed by loading
74 the contents of Supplementary File 1 with QIIME 2 View, providing far more detailed information than can
75 typically be provided in Methods text. QIIME 2 results are also semantically typed (Figure 2) and actions

76 indicate acceptable input types, clarifying the data that actions should be applied to and making complex workflows less error-prone.

78 Finally, QIIME 2 provides a software development kit (see <https://dev.qiime2.org>) that can be used to integrate
it as a component of other systems (e.g., such as Qiita²¹ or Illumina BaseSpace) and to develop interfaces
80 targeted toward users with different levels of computational sophistication (Figure S2). QIIME 2 provides the
QIIME 2 Studio graphical user interface and QIIME 2 View, interfaces designed for end-user biologists,
82 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; and q2cli and q2cwl, providing a
84 command line interface and Common Workflow Language²² wrappers for QIIME 2, designed for
high-performance computing experts.

86 There are many other powerful open source software tools for microbiome data science, including mothur²³,
phyloseq²⁴ and related tools available through Bioconductor²⁵, and the biobakery suite^{19,20,26}. mothur is a
88 microbiome bioinformatics platform that is often compared to QIIME 1 and QIIME 2. A major difference
between the two lies in the interactive visualizations: QIIME 2 provides many interactive visualization tools
90 (several examples are provided in Figure 1), while mothur focuses on generating data that can be easily
loaded and visualized with other tools. phyloseq focuses on microbiome statistical analysis and generating
92 publication-ready visualizations but, unlike QIIME 2, begins with a feature or OTU table, leaving "upstream"
processing steps such as sequence demultiplexing and quality control to other processing pipelines, many of
94 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
96 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
98 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
100 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 File 3.

102 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 Notebooks²⁷. For example, the BIOM
104 format²⁸ is supported by all of them. A diverse ecosystem of interoperable software is beneficial for the field, as
it allows experienced users to get multiple perspectives on their data and novice bioinformaticians to work in
106 programming environments that they are most comfortable with (e.g., phyloseq allows users to work in R, while
QIIME 2 allows users to work in Python). We plan to continue working with the developers of these tools and
108 organizations such as the Genomics Standards Consortium on plugins and standards to ensure
interoperability.

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

Figures and figure captions

114

116

118

120

122

124

126

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 Online Methods. (A) Unweighted UniFrac PCoA 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) A feature volatility plot illustrating change in *Bifidobacterium* abundance over time in breast-fed and formula-fed infants. Temporally interesting features can be interactively discovered with this visualization. (C) Interactive taxonomic composition bar plot illustrating 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 reduce the burden of exploratory analysis over QIIME 1. (D) Molecular cartography of the human skin surface. Colored spots represent the abundance of the small molecule cosmetic, sodium laureth sulfate, on the human skin. Sample data can be interactively visualized on 3D models, supporting the discovery of spatial patterns.

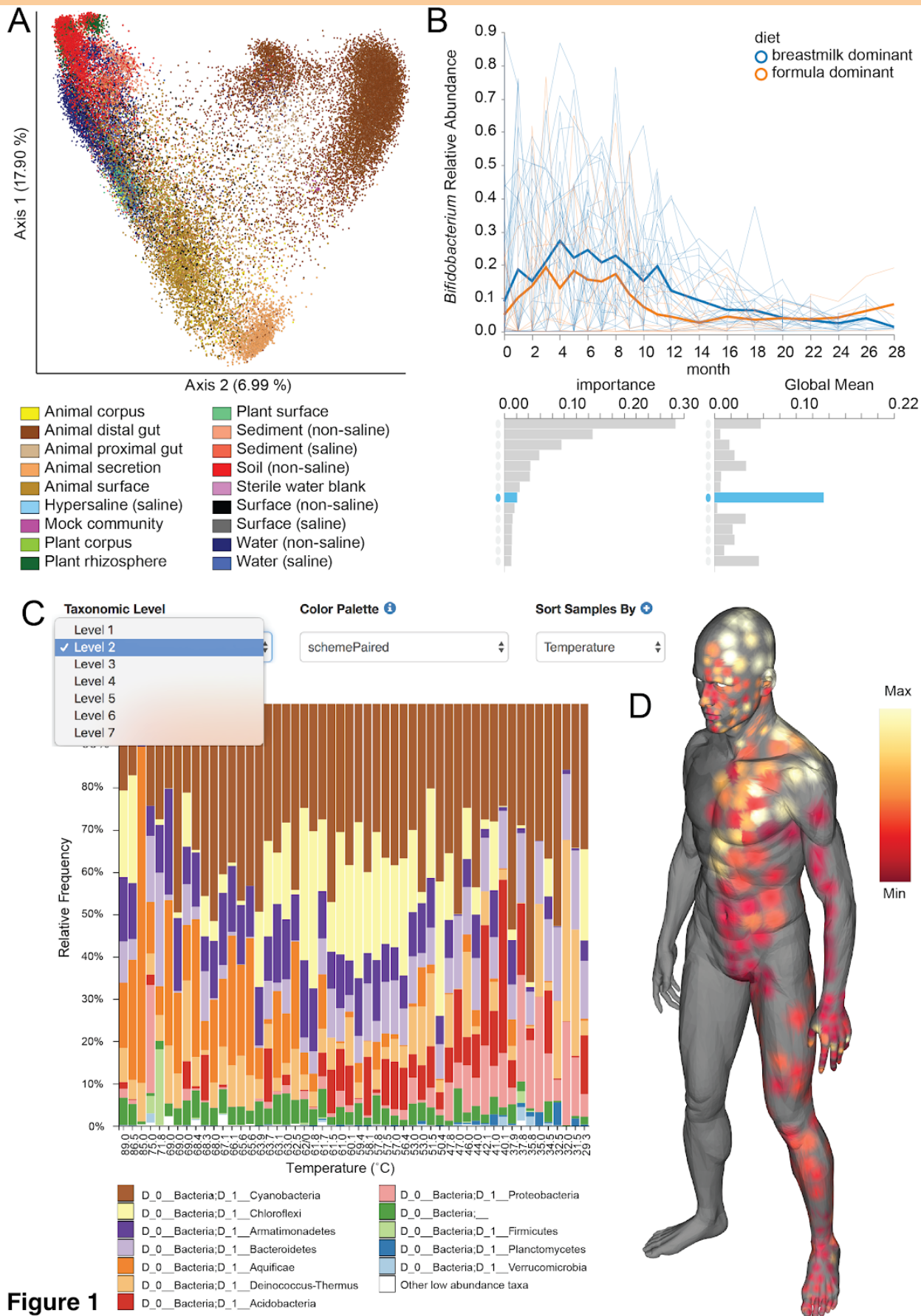


Figure 1

Figure 2: QIIME 2 iteratively records data provenance, ensuring bioinformatics reproducibility. This simplified diagram illustrates the automatically tracked information about the creation of the taxonomy barplot presented in Figure 1c. 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 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 as they were generated from the same upstream analysis. The resulting Visualization (Figure 1c), has the complete data provenance and correctly identifies shared processing of inputs. An interactive and complete version of this provenance graph (as well as those for other Figures 1 panels) can be accessed through Supplementary File 1.

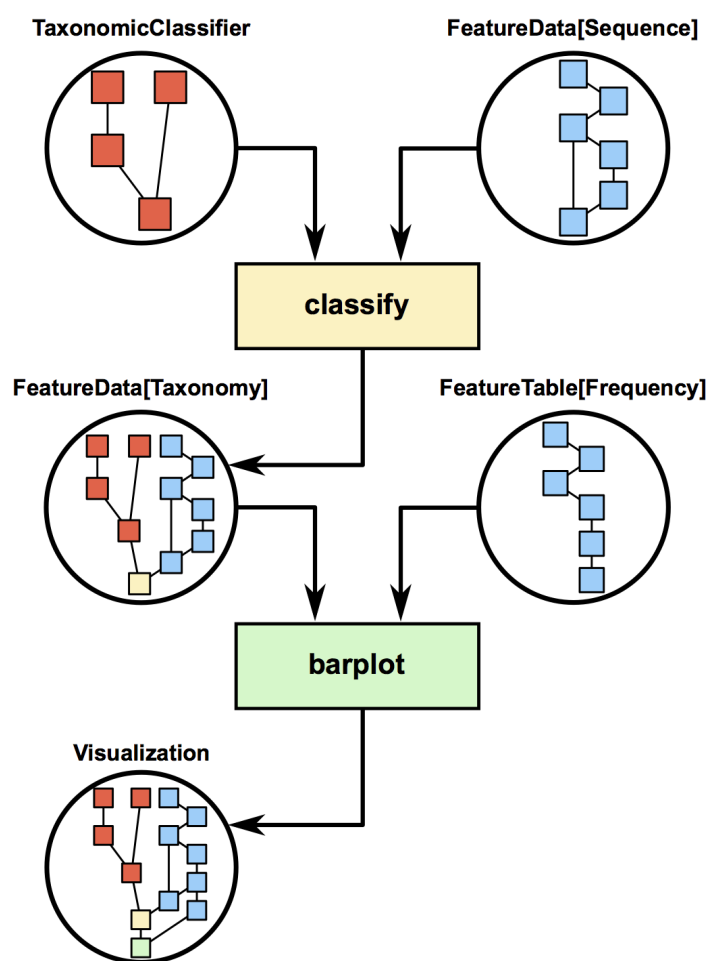


Figure 2

Code availability

142 QIIME 2 is open source and free for all use, including commercial. It is licensed under the BSD 3-clause
license. Source code is available at <https://github.com/qiime2>.

Data availability

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

Acknowledgements

152 QIIME 2 development was primarily funded by NSF Awards 1565100 to JGC and 1565057 to RK. Partial
support was also provided from the following grants: NIH U54CA143925 (JGC, TP) and U54MD012388 (JGC,
154 TP); grants from the Alfred P. Sloan Foundation (JGC, RK); ERC-STG project MetaPG (NS); Strategic Priority
Research Program of the Chinese Academy of Sciences QYZDB-SSW-SMC021 (YB); from the Australian
156 National Health and Medical Research Council APP1085372 (GAH, JGC, Von Bing Yap and RK); and from
Natural Sciences and Engineering Research Council (NSERC) to DLG; and under the State of Arizona
158 Technology and Research Initiative Fund (TRIF), administered by the Arizona Board of Regents, through
Northern Arizona University. All NCI co-authors were supported by the Intramural Research Program of the
160 National Cancer Institute. Thanks to the Yellowstone Center for Resources for research permit #5664 to JRS for
Yellowstone access and sample collection. We thank Paul J. McMurdie for helpful discussion on the
162 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 Ahmed Abdelfattah (Stockholm
164 University, Sweden), Rozlyn C.T. Boutin (University of British Columbia, Canada), David J. Bradshaw II (Florida
Atlantic University Harbor Branch Oceanographic Institute, USA), Lorinda Bullington (MPG Ranch, USA),
166 Justine W. Debelius (Karolinska Institutet, Sweden), Claire Duvallat (Massachusetts Institute of Technology,
USA), Erika Korzune Ganda (Cornell University, USA), Alexander Mahnert (Medical University of Graz,
168 Austria), Melanie C Melendrez (St. Cloud State University, USA), Devon O'rourke (University of New
Hampshire, USA), Adam R. Rivers (USDA-ARS, USA), Biswarup Sen (Tianjin University, China), Solveig
170 Tangedal (Haukeland University Hospital and University of Bergen, Norway), Pedro J. Torres (San Diego State
University, USA), and Jonathan Warren (National Laboratory Service, UK) for writing end-user reviews
172 included in Supplementary File 3.

Author contributions

174 EB, JRR, MRD, NAB, YB, JEB, CJB, AMC, EC, RD, CFE, MEs, JMG, DLG, AKJ, KBK, STK, IK, TK, JL, YL,
AVM, JLM, LFN, SBO, DP, AS, SJS, ADS, LRT, PJTo, PJTu, SU, FV, JW, RK, and JGC developed
176 documentation, educational materials, and/or user/developer support content. EB, JRR, MRD, NAB, RK, and

JGC wrote the manuscript; all authors assisted with revision of the manuscript. EB, JRR, MRD, NAB, and JGC designed and developed the QIIME 2 framework. DMD, RL, EL, SCM, RS, JRS, WW, CHDW, and RK contributed data used in the manuscript and/or testing of the QIIME 2. CCA, CTB, EC, PCD, SH, PK, EL, TP, RS, EV, YW, and RK contributed to the design of analytic methods. EB, JRR, MRD, NAB, GAA, HA, EJA, MA, FA, KB, AB, BJC, JC, GMD, CD, MEr, JF, AG, KG, JG, BH, HH, CH, GH, SJ, LJ, BK, CRK, DK, JK, MGIL, CL, MM, CM, BM, DM, LJM, JM, ATN, JAN, SLP, MLP, EP, LBR, AR, MSR, PR, NS, MS, PT, AT, JJJV, YV, MV, MW, KCW, ADW, ZZX, JRZ, YZ, QZ, and JGC contributed software to QIIME 2 plugins, interfaces, framework, and/or build and test systems.

Main text references

1. Smith, M.I. et al. *Science* **339**, 548–554 (2013).
2. Gopalakrishnan, V. et al. *Science* **359**, 97–103 (2018).
3. Gehring, C.A., Sthultz, C.M., Flores-Rentería, L., Whipple, A.V. & Whitham, T.G. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 11169–11174 (2017).
4. Lee, K., Pletcher, S.D., Lynch, S.V., Goldberg, A.N. & Cope, E.K. *Front. Cell. Infect. Microbiol.* **8**, 168 (2018).
5. Metcalf, J.L. et al. *Science* **351**, 158–162 (2016).
6. Rubin, R.L. et al. *Ecol. Appl.* **28**, 1594–1605 (2018).
7. Pineda, A., Kaplan, I. & Bezemer, T.M. *Trends Plant Sci.* **22**, 770–778 (2017).
8. Kapono, C.A. et al. *Sci. Rep.* **8**, 3669 (2018).
9. Barr, T. et al. *Gut Microbes* 1–44 (2018).
10. Callahan, B.J. et al. *Nat. Methods* (2016).doi:10.1038/nmeth.3869
11. Amir, A. et al. *mSystems* **2**, (2017).
12. Bokulich, N.A. et al. *Microbiome* **6**, 90 (2018).
13. Janssen, S. et al. *mSystems* **3**, e00021–18 (2018).
14. Bokulich, N.A. et al. *mSystems* **3**, e00219–18 (2018).
15. Bokulich, N. et al. *JOSS* **3**, 934 (2018).
16. Sedio, B.E., Rojas Echeverri, J.C., Boya P, C.A. & Wright, S.J. *Ecology* **98**, 616–623 (2017).
17. Wang, M. et al. *Nat. Biotechnol.* **34**, 828–837 (2016).
18. Hillmann, B. et al. *bioRxiv* 320986 (2018).doi:10.1101/320986
19. Truong, D.T. et al. *Nat. Methods* **12**, 902–903 (2015).
20. Langille, M.G.I. et al. *Nat. Biotechnol.* **31**, 814–821 (2013).
21. Gonzalez, A. et al. *Nat. Methods* **15**, 796–798 (2018).
22. Amstutz, P. et al. (2016).doi:10.6084/m9.figshare.3115156.v2
23. Schloss, P.D. et al. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009).
24. McMurdie, P.J. & Holmes, S. *PLoS One* **8**, e61217 (2013).
25. Huber, W. et al. *Nat. Methods* **12**, 115–121 (2015).
26. Franzosa, E.A. et al. *Nat. Methods* **15**, 962–968 (2018).
27. Kluiver, T. et al. *Positioning and Power in Academic Publishing: Players, Agents and Agendas* 87–90 (2016).
28. McDonald, D. et al. *Gigascience* **1**, 7 (2012).

Author affiliations

¹Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA. ²Metabolic Epidemiology Branch, National Cancer Institute, Rockville, MD, USA. ³Department of Computer Science and Engineering, University of Minnesota, Minneapolis, Minnesota, USA. ⁴Biology Department, Woods Hole Oceanographic Institution, Woods Hole, MA, USA. ⁵Department of Population Health and Reproduction, University of California, Davis, CA, USA. ⁶Department of Biological Engineering, Massachusetts Institute of Technology,

Cambridge, MA, USA. ⁷Center for Microbiome Informatics and Therapeutics, Massachusetts Institute of Technology, Cambridge, MA, USA. ⁸University of Copenhagen, Faculty of Health and Medical Sciences, Novo Nordisk Foundation Center for Basic Metabolic Research, Copenhagen, Denmark. ⁹Centre for Integrative Biology, University of Trento, Trento, Italy. ¹⁰State Key Laboratory of Plant Genomics, Institute of Genetics and Developmental Biology, Chinese Academy of Sciences, Beijing, China. ¹¹Centre of Excellence for Plant and Microbial Sciences (CEPAMS), Institute of Genetics and Developmental Biology, Chinese Academy of Sciences & John Innes Centre, Beijing, China. ¹²University of Chinese Academy of Sciences, Beijing, China. ¹³Department of Microbiology and Immunology, University of California, San Francisco, CA, USA. ¹⁴Division of Gastroenterology and Nutrition, Children's Hospital of Philadelphia, Philadelphia, PA, USA. ¹⁵Hepatology, Children's Hospital of Philadelphia, Philadelphia, PA, USA. ¹⁶Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Denmark. ¹⁷Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, WA, USA. ¹⁸Department of Population Health & Pathobiology, North Carolina State University, Raleigh, NC, USA. ¹⁹Bioinformatics Research Center, North Carolina State University, Raleigh, NC, USA. ²⁰Collaborative mass spectrometry innovation center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, San Diego, CA, USA. ²¹Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA. ²²Department of Microbiology and Immunology, Dalhousie University, Halifax, Nova Scotia, Canada. ²³Irving K. Barber School of Arts and Sciences, University of British Columbia, Kelowna, British Columbia, Canada. ²⁴A. Watson Armour III Center for Animal Health and Welfare, Aquarium Microbiome Project, John G. Shedd Aquarium, Chicago, IL, USA. ²⁵Department of Biology, University of British Columbia Okanagan, Okanagan, BC, Canada. ²⁶Computational Bioscience Graduate Program, University of Colorado Denver Anschutz Medical Campus, Aurora, Colorado, USA. ²⁷Department of Medicine, Division of Biomedical Informatics and Personalized Medicine, University of Colorado Denver Anschutz Medical Campus, Aurora, Colorado, USA. ²⁸Irving K. Barber School of Arts and Sciences, Department of Biology, The University of British Columbia, Kelowna, BC, Canada. ²⁹Department of Medicine, The University of British Columbia, Kelowna, BC, Canada. ³⁰Department of Pediatrics, University of California San Diego, La Jolla, CA, USA. ³¹Center for Microbial Ecology, Michigan State University, East Lansing, MI, USA. ³²Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA. ³³Stanford University, Statistics Department, Palo Alto, CA, USA. ³⁴Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA. ³⁵Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA. ³⁶Broad Institute of MIT and Harvard, Cambridge, MA, USA. ³⁷Research School of Biology, The Australian National University, Canberra, ACT, Australia. ³⁸Department of Pediatric Oncology, Hematology and Clinical Immunology, Heinrich-Heine University Dusseldorf, Dusseldorf, Germany. ³⁹Department of Family Medicine and Public Health, University of California San Diego, La Jolla, CA, USA. ⁴⁰College of Pharmacy, Sookmyung Women's University, Seoul, Republic of Korea. ⁴¹San Diego State University, Department of Biology, San Diego, CA, USA. ⁴²Biotechnology Institute, University of Minnesota, Saint Paul, MN, USA. ⁴³Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA. ⁴⁴Department of Pediatrics, University of California San Diego, La Jolla, CA, USA. ⁴⁵Department of Pharmacology, Dalhousie University, Halifax, Nova Scotia, Canada. ⁴⁶Science Education, Howard Hughes Medical Institute, Ashburn, VA, USA. ⁴⁷Department of Microbiome Science, Max Planck Institute for Developmental Biology, Tübingen, Germany. ⁴⁸Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY, USA. ⁴⁹Department of Medicine, Division of Biomedical Informatics and Personalized Medicine, University of Colorado Denver Anschutz Medical Campus, Aurora, CO, USA. ⁵⁰Department of Computer Science & Engineering, University of California San Diego, La Jolla, CA, USA. ⁵¹Department of Statistics, University of Washington, Seattle, WA, USA. ⁵²Department of Animal Science, Colorado State University, Fort Collins, CO, USA. ⁵³Irving K. Barber School of Arts and Sciences, Unit 2 (Biology), University of British Columbia, Kelowna, BC, Canada. ⁵⁴Mountain View, Google LLC, Mountain View, CA, USA. ⁵⁵Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, USA. ⁵⁶School of Information Studies, Syracuse University, Syracuse, NY, USA. ⁵⁷School of STEM, University of Washington Bothell, Bothell, WA, USA. ⁵⁸Department of Biological Sciences, Webster University, St Louis, MO, USA. ⁵⁹Agricultural Research Service, Genomics and Bioinformatics Research Unit, United States Department of Agriculture, Gainesville, FL, USA. ⁶⁰College of Medicine, Department of Biomedical Informatics, University of Arkansas for Medical Sciences, Little Rock, AR, USA. ⁶¹Computational Bioscience Program, University of Colorado Denver Anschutz Medical Campus, Aurora, CO, USA. ⁶²Department of Civil

276 and Environmental Engineering, Colorado School of Mines, Golden, CO, USA. ⁶³Department of Biological
Sciences and Northern Gulf Institute, University of Southern Mississippi, Hattiesburg, Mississippi, USA.
278 ⁶⁴Ocean Chemistry and Ecosystems Division, Atlantic Oceanographic and Meteorological Laboratory, National
Oceanic and Atmospheric Administration, La Jolla, CA, USA. ⁶⁵Department of Biology, San Diego State
280 University, San Diego, CA, USA. ⁶⁶Department of Environmental and Occupational Health Sciences, University
of Washington, Seattle, WA, USA. ⁶⁷Division of Biological Sciences, University of California San Diego, San
282 Diego, CA, USA. ⁶⁸Department of Microbiology and Immunology, University of California San Francisco, San
Francisco, CA, USA. ⁶⁹Quantitative and Systems Biology Graduate Program, University of California Merced,
284 Merced, CA, USA. ⁷⁰Bioinformatics Group, Wageningen University, Wageningen, The Netherlands.
⁷¹Department of Mathematics, University of Arizona, Tucson, AZ, USA. ⁷²National Laboratory Service,
286 Environment Agency, Starcross, UK. ⁷³College of Agriculture and Life Sciences, University of Florida,
Gainesville, FL, USA. ⁷⁴Department of Biostatistics, University of Washington, Seattle, WA, USA. ⁷⁵University of
288 Washington Bothell, School of STEM, Division of Biological Sciences, Bothell, WA, USA. ⁷⁶Merck & Co. Inc.,
Kenilworth, NJ, USA. ⁷⁷Department of Computer Science and Engineering, University of California San Diego,
290 La Jolla, California, USA.

Online Methods

2

Overview of QIIME 2

We provide a high-level overview of the QIIME 2 system. `Monospace font` is used to indicate literal terms, such as objects defined by QIIME 2. The most up-to-date information on these topics is available in the QIIME 2 developer documentation at <https://dev.qiime2.org>.

There are three core components of the QIIME 2 system architecture: the **framework**, the **interfaces**, and the **plugins** (Figure S1). **Interfaces** are responsible for turning user intent into action. **Plugins** define all domain-specific functionality. The most important restriction of the architecture, which is evident in Figure S1, is that interfaces and plugins do not communicate directly with one another -- that communication is always mediated by the **framework**. In other words, the domain-specific analytic functionality (defined in plugins) is entirely decoupled from how users interface with the system (defined in interfaces). This important constraint allows multiple kinds of interfaces to be dynamically generated, and as a result QIIME 2 can adapt its user interface to the audience and the task at hand (Figure S2).

Third-party developers can create and distribute both plugins and interfaces for QIIME 2 independently of the core QIIME 2 development group, which forms the basis for our goal of decentralized QIIME 2 development (see <https://library.qiime2.org> and <https://dev.qiime2.org>). By removing our team as a bottleneck in developers delivering their new methods to users through QIIME 2, microbiome research can advance more quickly by ensuring that QIIME 2 users can have access to the latest microbiome analytic methods as quickly as bioinformatics researchers and developers can distribute them. This model makes QIIME 2 (and tools that build on it, such as Qiita) a platform for microbiome data science, not only a tool for a specific type of analysis. Since plugins conform to requirements specified by the framework, framework features such as data provenance tracking and multiple interface support are available for all plugins without the plugin developers having to be aware of these features.

In the terminology of QIIME 2, an **Action** creates a **Result**, and a **Result** can be either an **Artifact** or a **Visualization**. An **Artifact** is data generated by one or more QIIME 2 **Actions** which can be used as the input to other QIIME 2 **Actions**. A **Visualization** on the other hand is a terminal output of QIIME 2, which could be an interactive visualization (as in the Figure 1 examples) or any other result that is intended to be consumed by humans (not by a QIIME 2 **Action**). QIIME 2 assigns version 4 universally unique identifiers (UUIDs) to each execution of an **Action**, and to all **Results**. QIIME 2 stores information about the series of **Actions** that led to a **Result**, along with information about the environment (including versions of all QIIME 2 packages and other Python dependencies) where each **Action** was executed, and the data itself. We refer to this process as data provenance tracking, or simply provenance tracking. We did not want to create new bioinformatics file formats to support the storage of data provenance, so QIIME 2 **Results** are instead stored as zip files containing a data directory that contains only the data in a relevant format (e.g., fasta or fastq for sequence data, newick for phylogenetic trees, etc), plus QIIME-2-specific metadata in other directories (such as provenance). These files use the extension `.qza` (for QIIME zipped artifact) or `.qzv` (for QIIME zipped visualization), but they are standard zip files that could be unzipped using common tools such as unzip, WinZip, or 7-Zip. Additional motivations for the storage of QIIME 2 **Results** in these structured zip files include the ability to submit as supplementary material to journals (the extension can simply be changed to `.zip` if required by the journal); “future-proofing” of QIIME 2 **Results** (even if QIIME 2 weren’t used anymore, **Results** could still be accessed by unzipping `.qza` or `.qzv` files - see Extracting data from QIIME 2 archives below); zip files contain an index, allowing them to be inspected for certain information without uncompressing them; and data are always compressed, facilitating data sharing. Because provenance is stored alongside data in `.qza` and `.qzv` files, provenance tracking is decentralized (no QIIME 2

server or database needs to be keeping track of this information) ensuring that information on how data was generated will not be lost as long as the data is intact. However, assignment of UUIDs to all QIIME 2 Results (as described above) lends itself to managing these data in a database if that is desired.

Another important component of QIIME 2 is its **semantic type system**. All Artifacts used in QIIME 2 are annotated with a semantic description of their type which conveys the meaning of the data. Semantic types differ from data types (how data is represented in memory) or file formats (how data is stored on disk), and allow QIIME 2 to constrain the composition of multiple actions to only those combinations which are semantically meaningful without needing to consider the specific file formats or data types. This also makes it possible to determine what Actions could be applied (and in what order) to generate a given Artifact from some set of input Artifacts. For example, phylogenetic trees in QIIME 2 can be either rooted or unrooted, and these two concepts are represented by the semantic types `Phylogeny[Rooted]` and `Phylogeny[Unrooted]`, respectively. QIIME 2 could support loading these into multiple different data types, including a scikit-bio `TreeNode` object or an ete3 `Tree` object. Both of these types are typically stored on disk in a newick-formatted file, but this format doesn't contain easily accessible information on whether the phylogeny is rooted or unrooted. Some QIIME 2 Actions can only generate a `Phylogeny[Unrooted]` (such as `fasttree`), and some other Actions only work on `Phylogeny[Rooted]` (such as `beta-phylogenetic`, which computes UniFrac distances). The semantic type system allows QIIME 2 to determine that the output of `fasttree` should not be directly provided as input to `beta-phylogenetic`, and to provide the user with that information prior to execution. This can help a researcher who is new to microbiome data science avoid using data incorrectly. This will also enable QIIME 2 to automatically assist users in identifying relevant workflows to generate desired data or further explore data they already have.

Due to recent advances in package management systems and bioinformatics package repositories (e.g., Anaconda, Bioconda¹, and Bioconductor²), QIIME 2 is straightforward to install.

QIIME 2 View

QIIME 2 View (<https://view.qiime2.org>) is a unique and novel contribution to the microbiome data science ecosystem that facilitates collaborative research. A user who has generated QIIME 2 visualizations can share those visualizations with a collaborator who can explore the results interactively without having QIIME 2 installed. QIIME 2 View achieves this simplified sharing of complex interactive visualizations through a novel combination of modern web browser APIs within a single-page application. It allows a user's browser to open and read `.qza` and `.qzv` files without the need to transfer the files over the network by utilizing a Service Worker to redirect HTTP requests directly into the archive which is retained on the user's computer. This approach of data unpackaging and local command execution makes QIIME 2 View well suited to cases where the results are unpublished or contain private information (that information will not be stored on any remote server). It is also possible to create "smart" URLs which automatically fetch content from a CORS-enabled web-server (for example, see the links in the README.md file at <https://github.com/qiime2/paper1>). This makes it very simple to share a single link with a collaborator that will be resolved into a fully interactive visualization on a user's computer automatically. The structured nature of the archive format (Figure S3) also allows QIIME 2 View to generate a dynamic provenance visualization, summarizing the entire provenance of the archive in question.

Extracting data from QIIME 2 archives

QIIME 2 `.qza` and `.qzv` files are zip file containers with a defined internal directory structure. It's very easy to get data out in the canonical formats (Figure S3). If QIIME 2 and the `q2cli` command line interface are installed,

this can be achieved using the `qiime tools export` command. If QIIME 2 is not installed, this can be achieved using standard decompression utilities such as `unzip`, WinZip, or 7-zip. We illustrate how this can be achieved using `unzip` on macOS. This can similarly be achieved on Windows or Linux. We illustrate this here to further future-proof QIIME 2 Results - even if the QIIME 2 documentation were no longer accessible, users could follow these steps to access QIIME 2 Results.

First, obtain a `.qza` file. Here we use the `FeatureData[Sequence]` artifact generated during the QIIME 2 Moving Pictures tutorial.

```
$ wget https://docs.qiime2.org/2018.8/data/tutorials/moving-pictures/rep-seqs.qza
```

Next, `unzip` that file with the macOS (or Linux) `unzip` program. This will create a new directory. The name of that directory will be the UUID of the artifact being unzipped, in this case

```
8dc793b8-7284-462a-8578-6370ffccebd
```

```
$ unzip rep-seqs.qza
```

```
Archive:  rep-seqs.qza
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/metadata.yaml
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/VERSION
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/provenance/metadata.yaml
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/provenance/citations.bib
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/provenance/VERSION
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/metadata.yaml
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/citations.bib
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/VERSION
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/action/action.yaml
```

```
l
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/metadata.yaml
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/citations.bib
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/VERSION
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/action/action.yaml
```

```
l
```

```
  inflating:
```

```
8dc793b8-7284-462a-8578-6370ffccebd/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/action/barcodes.tsv
```

```
sv
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/provenance/action/action.yaml
```

```
  inflating: 8dc793b8-7284-462a-8578-6370ffccebd/data/dna-sequences.fasta
```

The last entry that is unzipped in this example is `data/dna-sequences.fasta`. All other directories and files are QIIME 2 specific metadata (such as information about the semantic type of the artifact and the data provenance). If you're only interested in the sequence data, you can safely ignore all of that information. The `data/dna-sequences.fasta` file is a typical fasta file containing sequence identifiers and sequences. The first four lines of this file can be viewed as follows:

```
$ head -4 8dc793b8-7284-462a-8578-6370ffccebd/data/dna-sequences.fasta
```

```
>f352c1f1efecf483511c2270aab0ae6
```

```
TACGTAGGGTGCAGCGTTAATCGGAATTACTGGCGTAAAGCGTGCAGCGGTTTTGTAAGACAGAGGTGAAATCCCCGGGCTCAACCTGGGAACCTG
```

```
CCTTTGTGACTGCAAGGCTG
```

```
>82e72255267397b777a1afd44ea22755
```

138 TACGGAGGATCCAAGCGTTATCCGGAATCATTGGGTTTAAAGGGTCCGTAGGCGGTTTAGTAAGTCAGTGGTGAAAGCCATCGCTCAACGGTGGAAACGG
140 CCATTGATACTGCTAGACTT

142 QIIME 2 user and developer community

144 QIIME 2 officially succeeded QIIME 1 (<http://www.qiime.org>)³ in January of 2018, and has developed an
146 engaged user base and community. As of this writing there are over 1980 active users (users who have
148 performed an action, such as creating or liking a post) on the QIIME 2 Forum; over 3000 monthly downloads of
QIIME 2 from Anaconda; over 8000 unique visitors to the QIIME 2 Forum according to Google Analytics; and
our multi-day workshops are frequently filled to capacity (<https://workshops.qiime2.org>). QIIME 2 is also being
adopted by third-party bioinformatics developers who are choosing to make their software accessible through
plugins, and who are motivated to develop for QIIME 2 by access to its integrated provenance tracking,
multiple interfaces, standardization of data types provided by the semantic type system, large user community,
and supportive developer community.

150 A core goal of QIIME 2 is to cultivate a diverse and inclusive community of scientists, software engineers,
152 statisticians, educators, students, and other microbiome stakeholders who are openly sharing methods, data,
and knowledge to advance microbiome research.

Supplementary figures, files, and captions

154 **Figure S1. Schematic diagram of the QIIME 2 system.** Interfaces define how users interact with the system;
156 plugins define all domain-specific functionality; and the framework mediates communication between plugins
158 and interfaces, and performs core functionality such as provenance tracking. Arrows indicate dependencies.
Interfaces interact only with the `qiime2.sdk` submodule, while plugins interact only with the `qiime2.plugin`
submodule. This design has led to a system that is readily extended by third-party plugin and interface
developers.

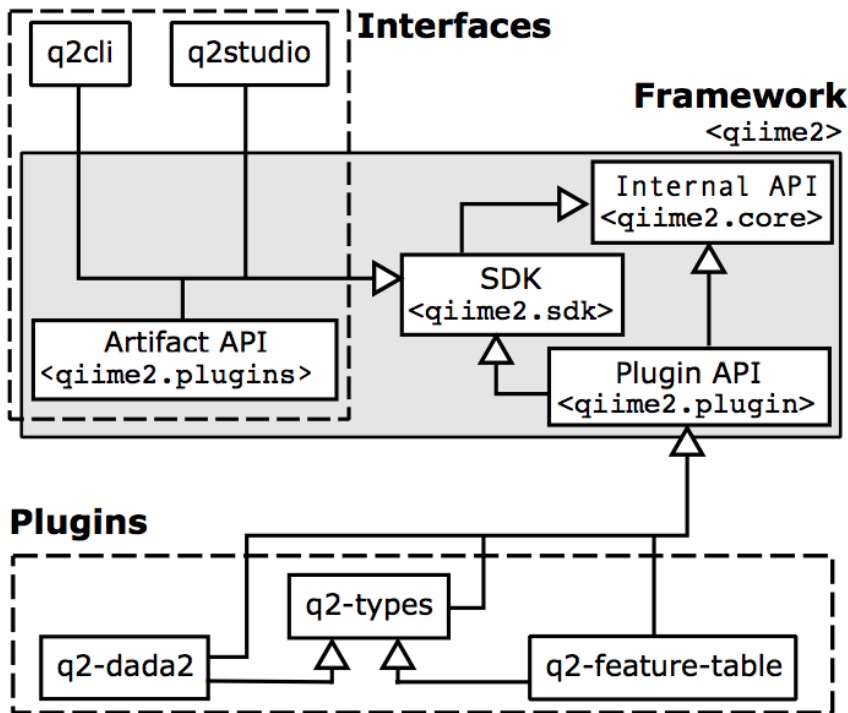
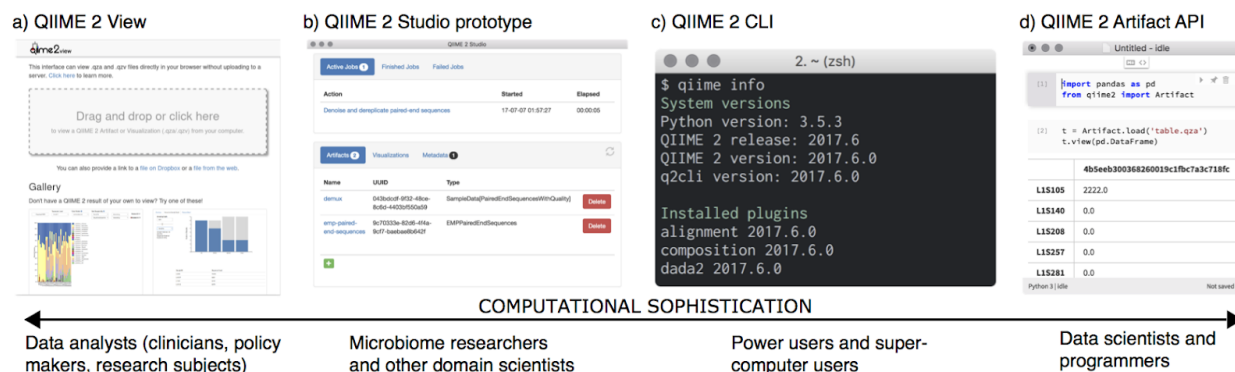


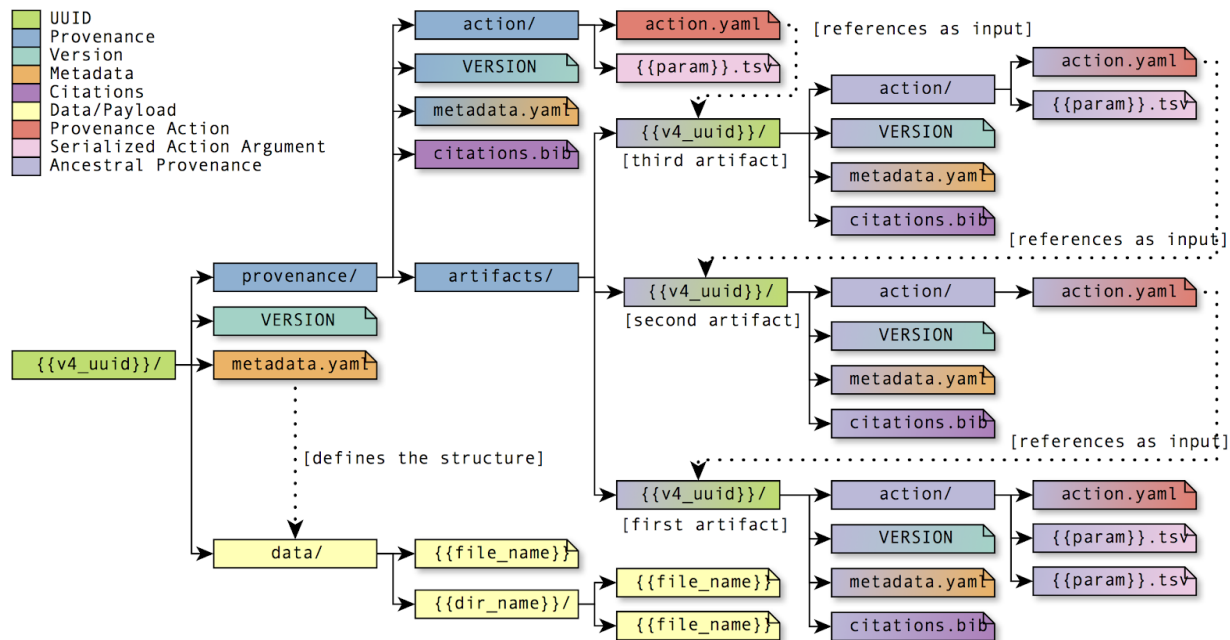
Figure S2. QIIME 2 is interface agnostic. The full suite of QIIME 2 functionality is useful to and usable by researchers ranging widely in their computational sophistication, a major advantage over technologies such as QIIME 1 that provide a single interface. (a) Users wanting to view QIIME 2 results or data provenance can use QIIME 2 View without installing QIIME 2, which is convenient for lead investigators, clinicians, or policy makers who may want to explore interactive visualizations generated by others. (b) Researchers who prefer graphical interfaces can use QIIME 2 Studio, our prototype graphical interface. This is convenient for users without command line or programming skills. (c) Power users (e.g., who are comfortable with the Linux command line and/or regularly work on institutional computer clusters), can use QIIME 2 through the command line interface, q2cli. (d) “Data scientists” (e.g., users who are programmers, who work in Jupyter Notebooks, or who are interested in automating QIIME 2 workflows), can use QIIME 2 through the Python 3 “artifact API”.



170

Figure S3. Anatomy of a QIIME 2 Archive (i.e., .qza or .qzv file). QIIME 2 stores data in a directory structure called an Archive. These archives are zipped to make moving data convenient. The directory structure has a single root directory named with a UUID which serves as the identity of the archive.

172



174 **Supplementary File 1** contains the QIIME 2 .qzv files corresponding to **Figure 1a-d**.

Supplementary File 2 contains four worked examples of using QIIME 2.

Supplementary File 3 contains notes from end users on QIIME 2.

176

Online methods references

- 178 1. Grünig, B. et al. *Nat. Methods* **15**, 475–476 (2018).
2. Huber, W. et al. *Nat. Methods* **12**, 115–121 (2015).
3. Caporaso, J.G. et al. *Nat. Methods* **7**, 335–336 (2010).
- 180 4. Thompson, L.R. et al. *Nature* **551**, 457–463 (2017).
5. McDonald, D. et al. *mSystems* **3**, e00031–18 (2018).
- 182 6. McDonald, D. et al. *Gigascience* **1**, 7 (2012).
7. Amir, A. et al. *mSystems* **2**, (2017).
- 184 8. Lozupone, C. & Knight, R. *Appl. Environ. Microbiol.* **71**, 8228–8235 (2005).
9. Vázquez-Baeza, Y., Pirrung, M., Gonzalez, A. & Knight, R. *Gigascience* **2**, 16 (2013).
- 186 10. Bokulich, N.A. et al. *Sci. Transl. Med.* **8**, 343ra82 (2016).
11. Callahan, B.J. et al. *Nat. Methods* (2016).doi:10.1038/nmeth.3869
- 188 12. McDonald, D. et al. *ISME J.* **6**, 610–618 (2012).
13. Bokulich, N.A. et al. *Microbiome* **6**, 90 (2018).
- 190 14. Bokulich, N. et al. *bioRxiv* 223974 (2017).doi:10.1101/223974
15. Protsyuk, I. et al. *Nat. Protoc.* **13**, 134–154 (2018).
- 192 16. Bouslimani, A. et al. *Proc. Natl. Acad. Sci. U. S. A.* **112**, E2120–9 (2015).

Supplementary File 2: QIIME 2 worked examples

Figure 1 presents the output of four worked examples of QIIME 2, and Supplementary File 1 contains the QIIME 2 `.qzv` files corresponding to Figure 1a-d. These are also accessible at <https://github.com/qiime2/paper1> and can be viewed using QIIME 2 View (<https://view.qiime2.org>) where readers can interact with the results, and explore the methods used to generate them (see the Provenance tab after loading a `.qzv` file with QIIME 2 View).

Here we describe the methods used to generate each of these visualizations at a high level, and present the QIIME 2 steps to generate these visualizations both as a series of command line interface (CLI) commands and as application programmer interface (API) calls. The steps presented are derived from the data provenance of each of the visualizations, and can be compared directly to the Provenance tab on QIIME 2 View. These worked examples represent real-world working conditions where different analysis steps were conducted at different times by different individuals running multiple releases of QIIME 2. This is common, for example as in Figure 2, when a researcher uses a pre-trained taxonomic classifier such as those available on the QIIME 2 website and forum. Because detailed information on system and plugin versions is tracked in provenance, consumers of QIIME 2 results and data provenance have comprehensive information about the software and dependency environment where results were generated ensuring reproducibility. In some examples data processed with other pipelines are imported into QIIME 2 (Figures 1a and 1d), while in others analysis begins with DNA sequencing data (Figures 1b and 1c). This illustrates that QIIME 2 can be used as an “end-to-end” microbiome analysis pipeline, or as a component of workflows that use other bioinformatics tools.

Figure 1a (a-pcoa.qzv in Supplementary File 1)

Emperor PCoA plot presenting a meta-analysis of the first release of the Earth Microbiome Project (EMP)⁴ and the first release of the American Gut Project (AGP)⁵. The EMP data was obtained from <ftp://ftp.microbio.me/emp/release1>, and the AGP data was obtained from Qiita study 10317 for the set of samples used in its publication (samples described in the AGP supplemental data accession table). Both projects were downloaded and imported into QIIME 2 as BIOM tables⁶. Those BIOM tables were combined, filtered for blooms⁷, rarefied at an even depth (1000 sequences per sample), and compared using the unweighted UniFrac⁸ metric. Lastly the samples were projected into a small dimensional space using principal coordinates analysis and visualized using Emperor⁹. The samples were colored according to the Earth Microbiome Project Ontology⁴.

Command line interface steps

```
# versions: {'types': '2018.8.0'}
qiime tools import \
  --type 'FeatureData[Sequence]' \
  --input-path combined-features.fasta \
  --input-format DNAFASTAFormat \
  --output-path combined-features.fasta.qza

# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}
qiime feature-table filter-seqs \
  --i-data combined-features.fasta.qza \
```

```
42 --m-metadata-file metadata1.txt \  
--p-exclude-ids \  
--output-dir feature-table-filter_seqs_1  
  
44 # versions: {'types': '2018.8.0', 'fragment-insertion': '2018.6.17'}  
qiime fragment-insertion sepp \  
46 --i-representative-sequences feature-table-filter_seqs_1/filtered_data.qza \  
--p-threads 20 \  
48 --p-alignment-subset-size 1000 \  
--p-placement-subset-size 5000 \  
50 --p-no-debug \  
--output-dir fragment-insertion-sepp_1  
  
52 # versions: {'types': '2018.8.0'}  
qiime tools import \  
54 --type 'FeatureTable[Frequency]' \  
--input-path emp.biom \  
56 --input-format BIOMV210Format \  
--output-path emp.biom.qza  
  
58 # versions: {'types': '2018.8.0'}  
qiime tools import \  
60 --type 'FeatureTable[Frequency]' \  
--input-path agp.upper.biom \  
62 --input-format BIOMV210Format \  
--output-path agp.upper.biom.qza  
  
64 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
qiime feature-table merge \  
66 --i-tables agp.upper.biom.qza \  
--i-tables emp.biom.qza \  
68 --p-overlap-method 'error_on_overlapping_sample' \  
--output-dir feature-table-merge_1  
  
70 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
qiime feature-table filter-features \  
72 --i-table feature-table-merge_1/merged_table.qza \  
--m-metadata-file metadata2.txt \  
74 --p-min-frequency 0 \  
--p-min-samples 0 \  
76 --p-exclude-ids \  
--output-dir feature-table-filter_features_1  
  
78 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
qiime feature-table rarefy \  
80 --i-table feature-table-filter_features_1/filtered_table.qza \  
--p-sampling-depth 1000 \  
82 --output-dir feature-table-rarefy_1
```

```
# versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
84 qiime diversity beta-phylogenetic-alt \
  --i-table feature-table-rarefy_1/rarefied_table.qza \
86  --i-phylogeny fragment-insertion-sepp_1/tree.qza \
  --p-metric 'unweighted_unifrac' \
88  --p-n-jobs 16 \
  --p-no-variance-adjusted \
90  --p-bypass-tips \
  --output-dir diversity-beta_phylogenetic_alt_1

92 # versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
qiime diversity filter-distance-matrix \
94  --i-distance-matrix diversity-beta_phylogenetic_alt_1/distance_matrix.qza \
  --m-metadata-file metadata3.txt \
96  --p-where 'project!="Not Available"' \
  --p-no-exclude-ids \
98  --output-dir diversity-filter_distance_matrix_1

# versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
100 qiime diversity pcoa \
  --i-distance-matrix
102 diversity-filter_distance_matrix_1/filtered_distance_matrix.qza \
  --p-number-of-dimensions 5 \
104  --output-dir diversity-pcoa_1

# versions: {'types': '2018.8.0', 'emperor': '2018.8.0'}
106 qiime emperor plot \
  --i-pcoa diversity-pcoa_1/pcoa.qza \
108  --m-metadata-file metadata4.txt \
  --output-dir emperor-plot_1
110
```

Application programmer interface steps

```
import qiime2
112 from qiime2.plugins import fragment_insertion
from qiime2.plugins import feature_table
114 from qiime2.plugins import diversity
from qiime2.plugins import emperor

116 metadata1_txt = qiime2.Metadata.load('metadata1.txt')
metadata2_txt = qiime2.Metadata.load('metadata2.txt')
118 metadata3_txt = qiime2.Metadata.load('metadata3.txt')
metadata4_txt = qiime2.Metadata.load('metadata4.txt')

120 # versions: {'types': '2018.8.0'}
combined_features_fasta = qiime2.Artifact.import_data(
122   'FeatureData[Sequence]', 'combined-features.fasta', view_type='DNAFASTAFormat'
```

```
)  
124 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
feature_table_filter_seqs_1 = feature_table.actions.filter_seqs(  
126     data=combined-features_fasta,  
     metadata=metadata1_txt,  
128     where=None,  
     exclude_ids=True,  
130 )  
  
# versions: {'types': '2018.8.0', 'fragment-insertion': '2018.6.17'}  
132 fragment_insertion_sepp_1 = fragment_insertion.actions.sepp(  
     representative_sequences=feature_table_filter_seqs_1.filtered_data,  
134     threads=20,  
     alignment_subset_size=1000,  
136     placement_subset_size=5000,  
     debug=False,  
138 )  
  
# versions: {'types': '2018.8.0'}  
140 emp_biom = qiime2.Artifact.import_data(  
     'FeatureTable[Frequency]', 'emp.biom', view_type='BIOMV210Format'  
142 )  
  
# versions: {'types': '2018.8.0'}  
144 agp_upper_biom = qiime2.Artifact.import_data(  
     'FeatureTable[Frequency]', 'agp.upper.biom', view_type='BIOMV210Format'  
146 )  
  
# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
148 feature_table_merge_1 = feature_table.actions.merge(  
     tables=agp_upper_biom,  
150     tables=emp_biom,  
     overlap_method='error_on_overlapping_sample',  
152 )  
  
# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
154 feature_table_filter_features_1 = feature_table.actions.filter_features(  
     table=feature_table_merge_1.merged_table,  
156     metadata=metadata2_txt,  
     min_frequency=0,  
158     max_frequency=None,  
     min_samples=0,  
160     max_samples=None,  
     where=None,  
162     exclude_ids=True,  
    )
```

```
164 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}
feature_table_rarefy_1 = feature_table.actions.rarefy(
166     table=feature_table_filter_features_1.filtered_table,
        sampling_depth=1000,
168 )

# versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
170 diversity_beta_phylogenetic_alt_1 = diversity.actions.beta_phylogenetic_alt(
        table=feature_table_rarefy_1.rarefied_table,
172     phylogeny=fragment_insertion_sepp_1.tree,
        metric='unweighted_unifrac',
174     n_jobs=16,
        variance_adjusted=False,
176     alpha=None,
        bypass_tips=True,
178 )

# versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
180 diversity_filter_distance_matrix_1 = diversity.actions.filter_distance_matrix(
        distance_matrix=diversity_beta_phylogenetic_alt_1.distance_matrix,
182     metadata=metadata3_txt,
        where='project!="Not Available"',
184     exclude_ids=False,
        )

186 # versions: {'types': '2018.8.0', 'diversity': '2018.8.0'}
diversity_pcoa_1 = diversity.actions.pcoa(
188     distance_matrix=diversity_filter_distance_matrix_1.filtered_distance_matrix,
        number_of_dimensions=5,
190 )

# versions: {'types': '2018.8.0', 'emperor': '2018.8.0'}
192 emperor_plot_1 = emperor.actions.plot(
        pcoa=diversity_pcoa_1.pcoa,
194     metadata=metadata4_txt,
        custom_axes=None,
196 )
```

Figure 1b (b-feature-volatility.qzv in Supplementary File 1)

198 Data were generated on five sequencing runs of V4 16S rRNA gene amplicons from the ECAM study¹⁰.
Forward reads were imported separately in EMPSingleEndDirFmt format, demultiplexed with q2-demux's
200 emp_single method, and denoised using q2-dada2's denoise_single method (trunc_len=150, other
parameters used default values)¹¹. Denoised feature tables and sequences were merged using
202 q2-feature-table's merge and merge-seqs methods, respectively. q2-feature-table's filter-samples
method was used to remove samples with fewer than 2000 sequences, and to perform metadata-based

204 filtering to retain only children's samples. A naive Bayes taxonomy classifier was trained on the Greengenes¹²
reference sequences (clustered at 99% similarity) using q2-feature-classifier's
206 fit-classifier-naive-bayes method¹³. This classifier was used to taxonomically classify the ECAM
ASVs using q2-feature-classifier's classify-sklearn method¹³. ASVs were collapsed based on genus-level
208 taxonomy using q2-taxa's collapse method. Temporally predictive features were identified using
q2-longitudinal's feature-volatility pipeline¹⁴ using default parameters. Data contained in this artifact
210 have been described in a previous publication¹⁴.

Command line interface steps

```
212 # versions: {'framework': '2017.2.0'}
qiime tools import \
214   --type 'FeatureData[Taxonomy]' \
   --input-path 99_otu_taxonomy.txt \
216   --input-format TaxonomyFormat \
   --output-path 99_otu_taxonomy.txt.qza

218 # versions: {'framework': '2017.2.0'}
qiime tools import \
220   --type 'FeatureData[Sequence]' \
   --input-path 99_otus.fasta \
222   --input-format DNAFASTAFormat \
   --output-path 99_otus.fasta.qza

224 # versions: {'feature-classifier': '2017.2.0'}
qiime feature-classifier fit-classifier-naive-bayes \
226   --i-reference-reads 99_otus.fasta.qza \
   --i-reference-taxonomy 99_otu_taxonomy.txt.qza \
228   --p-classify--alpha 0.01 \
   --p-classify--chunk-size -1 \
230   --p-classify--class-prior 'null' \
   --p-classify--fit-prior \
232   --p-feat-ext--analyzer 'char_wb' \
   --p-no-feat-ext--binary \
234   --p-feat-ext--decode-error 'strict' \
   --p-feat-ext--encoding 'utf-8' \
236   --p-feat-ext--input 'content' \
   --p-feat-ext--lowercase \
238   --p-feat-ext--n-features 8192 \
   --p-feat-ext--ngram-range '[8, 8]' \
240   --p-feat-ext--non-negative \
   --p-feat-ext--norm 'l2' \
242   --p-feat-ext--preprocessor 'null' \
   --p-feat-ext--stop-words 'null' \
244   --p-feat-ext--strip-accent 'null' \
   --p-feat-ext--token-pattern '(?u)\b\w\w+\b' \
246   --p-feat-ext--tokenizer 'null' \
   --output-dir feature-classifier-fit_classifier_naive_bayes_1
```



```
248 # versions: {'framework': '2017.4.0'}
qiime tools import \
250   --type 'RawSequences' \
   --input-path EMPSingleEndDirFmtimport_dir \
252   --input-format EMPSingleEndDirFmt \
   --output-path EMPSingleEndDirFmtimport_dir.qza

254 # versions: {'demux': '2017.4.0'}
qiime demux emp-single \
256   --i-seqs EMPSingleEndDirFmtimport_dir.qza \
   --m-barcodes-file metadata1.txt \
258   --m-barcodes-column 'ColumnName' \
   --p-no-rev-comp-barcodes \
260   --p-rev-comp-mapping-barcodes \
   --output-dir demux-emp_single_1

262 # versions: {'dada2': '2017.4.0'}
qiime dada2 denoise-single \
264   --i-demultiplexed-seqs demux-emp_single_1/per_sample_sequences.qza \
   --p-trunc-len 150 \
266   --p-trim-left 0 \
   --p-max-ee 2.0 \
268   --p-trunc-q 2 \
   --p-chimera-method 'pooled' \
270   --p-min-fold-parent-over-abundance 1.0 \
   --p-n-threads 1 \
272   --p-n-reads-learn 1000000 \
   --p-hashed-feature-ids \
274   --output-dir dada2-denoise_single_1

# versions: {'framework': '2017.4.0'}
276 qiime tools import \
   --type 'RawSequences' \
278   --input-path EMPSingleEndDirFmtimport_dir \
   --input-format EMPSingleEndDirFmt \
280   --output-path EMPSingleEndDirFmtimport_dir.qza

# versions: {'demux': '2017.4.0'}
282 qiime demux emp-single \
   --i-seqs EMPSingleEndDirFmtimport_dir.qza \
284   --m-barcodes-file metadata2.txt \
   --m-barcodes-column 'ColumnName' \
286   --p-no-rev-comp-barcodes \
   --p-rev-comp-mapping-barcodes \
288   --output-dir demux-emp_single_2

# versions: {'dada2': '2017.4.0'}
```

```
290 qiime dada2 denoise-single \  
    --i-demultiplexed-seqs demux-emp_single_2/per_sample_sequences.qza \  
292    --p-trunc-len 150 \  
    --p-trim-left 0 \  
294    --p-max-ee 2.0 \  
    --p-trunc-q 2 \  
296    --p-chimera-method 'pooled' \  
    --p-min-fold-parent-over-abundance 1.0 \  
298    --p-n-threads 1 \  
    --p-n-reads-learn 1000000 \  
300    --p-hashed-feature-ids \  
    --output-dir dada2-denoise_single_2  
  
302 # versions: {'framework': '2017.4.0'}  
qiime tools import \  
304    --type 'RawSequences' \  
    --input-path EMPSingleEndDirFmtimport_dir \  
306    --input-format EMPSingleEndDirFmt \  
    --output-path EMPSingleEndDirFmtimport_dir.qza  
  
308 # versions: {'demux': '2017.4.0'}  
qiime demux emp-single \  
310    --i-seqs EMPSingleEndDirFmtimport_dir.qza \  
    --m-barcodes-file metadata3.txt \  
312    --m-barcodes-column 'ColumnName' \  
    --p-no-rev-comp-barcodes \  
314    --p-rev-comp-mapping-barcodes \  
    --output-dir demux-emp_single_3  
  
316 # versions: {'dada2': '2017.4.0'}  
qiime dada2 denoise-single \  
318    --i-demultiplexed-seqs demux-emp_single_3/per_sample_sequences.qza \  
    --p-trunc-len 150 \  
320    --p-trim-left 0 \  
    --p-max-ee 2.0 \  
322    --p-trunc-q 2 \  
    --p-chimera-method 'pooled' \  
324    --p-min-fold-parent-over-abundance 1.0 \  
    --p-n-threads 1 \  
326    --p-n-reads-learn 1000000 \  
    --p-hashed-feature-ids \  
328    --output-dir dada2-denoise_single_3  
  
# versions: {'framework': '2017.4.0'}  
330 qiime tools import \  
    --type 'RawSequences' \  
332    --input-path EMPSingleEndDirFmtimport_dir \  
    --input-format EMPSingleEndDirFmt \  

```

```
334 --output-path EMPSingleEndDirFmtimport_dir.qza

# versions: {'demux': '2017.4.0'}
336 qiime demux emp-single \
  --i-seqs EMPSingleEndDirFmtimport_dir.qza \
338 --m-barcodes-file metadata4.txt \
  --m-barcodes-column 'ColumnName' \
340 --p-no-rev-comp-barcodes \
  --p-rev-comp-mapping-barcodes \
342 --output-dir demux-emp_single_4

# versions: {'dada2': '2017.4.0'}
344 qiime dada2 denoise-single \
  --i-demultiplexed-seqs demux-emp_single_4/per_sample_sequences.qza \
346 --p-trunc-len 150 \
  --p-trim-left 0 \
348 --p-max-ee 2.0 \
  --p-trunc-q 2 \
350 --p-chimera-method 'pooled' \
  --p-min-fold-parent-over-abundance 1.0 \
352 --p-n-threads 1 \
  --p-n-reads-learn 1000000 \
354 --p-hashed-feature-ids \
  --output-dir dada2-denoise_single_4

# versions: {'framework': '2017.4.0'}
356 qiime tools import \
358 --type 'RawSequences' \
  --input-path EMPSingleEndDirFmtimport_dir \
360 --input-format EMPSingleEndDirFmt \
  --output-path EMPSingleEndDirFmtimport_dir.qza

# versions: {'demux': '2017.4.0'}
362 qiime demux emp-single \
364 --i-seqs EMPSingleEndDirFmtimport_dir.qza \
  --m-barcodes-file metadata5.txt \
366 --m-barcodes-column 'ColumnName' \
  --p-no-rev-comp-barcodes \
368 --p-rev-comp-mapping-barcodes \
  --output-dir demux-emp_single_5

# versions: {'dada2': '2017.4.0'}
370 qiime dada2 denoise-single \
372 --i-demultiplexed-seqs demux-emp_single_5/per_sample_sequences.qza \
  --p-trunc-len 150 \
374 --p-trim-left 0 \
  --p-max-ee 2.0 \
376 --p-trunc-q 2 \
```

```
378 --p-chimera-method 'pooled' \  
--p-min-fold-parent-over-abundance 1.0 \  
--p-n-threads 1 \  
380 --p-n-reads-learn 1000000 \  
--p-hashed-feature-ids \  
382 --output-dir dada2-denoise_single_5  
  
# versions: {'feature-table': '2017.4.0'}  
384 qiime feature-table merge-seq-data \  
--i-data1 dada2-denoise_single_5/representative_sequences.qza \  
386 --i-data2 dada2-denoise_single_4/representative_sequences.qza \  
--output-dir feature-table-merge_seq_data_1  
  
# versions: {'feature-table': '2017.4.0'}  
388 qiime feature-table merge-seq-data \  
390 --i-data1 feature-table-merge_seq_data_1/merged_data.qza \  
--i-data2 dada2-denoise_single_3/representative_sequences.qza \  
392 --output-dir feature-table-merge_seq_data_2  
  
# versions: {'feature-table': '2017.4.0'}  
394 qiime feature-table merge-seq-data \  
--i-data1 feature-table-merge_seq_data_2/merged_data.qza \  
396 --i-data2 dada2-denoise_single_2/representative_sequences.qza \  
--output-dir feature-table-merge_seq_data_3  
  
# versions: {'feature-table': '2017.4.0'}  
398 qiime feature-table merge-seq-data \  
400 --i-data1 feature-table-merge_seq_data_3/merged_data.qza \  
--i-data2 dada2-denoise_single_1/representative_sequences.qza \  
402 --output-dir feature-table-merge_seq_data_4  
  
# versions: {'feature-classifier': '2017.5.0'}  
404 qiime feature-classifier classify-sklearn \  
--i-reads feature-table-merge_seq_data_4/merged_data.qza \  
406 --i-classifier feature-classifier-fit_classifier_naive_bayes_1/classifier.qza \  
--p-chunk-size 262144 \  
408 --p-n-jobs 4 \  
--p-pre-dispatch '2*n_jobs' \  
410 --p-confidence 0.7 \  
--output-dir feature-classifier-classify_sklearn_1  
  
# versions: {'feature-table': '2017.4.0'}  
412 qiime feature-table merge \  
414 --i-table1 dada2-denoise_single_5/table.qza \  
--i-table2 dada2-denoise_single_4/table.qza \  
416 --output-dir feature-table-merge_1  
  
# versions: {'feature-table': '2017.4.0'}
```

```
418 qiime feature-table merge \  
    --i-table1 feature-table-merge_1/merged_table.qza \  
420    --i-table2 dada2-denoise_single_3/table.qza \  
    --output-dir feature-table-merge_2  
  
422 # versions: {'feature-table': '2017.4.0'}  
qiime feature-table merge \  
424    --i-table1 feature-table-merge_2/merged_table.qza \  
    --i-table2 dada2-denoise_single_2/table.qza \  
426    --output-dir feature-table-merge_3  
  
# versions: {'feature-table': '2017.4.0'}  
428 qiime feature-table merge \  
    --i-table1 feature-table-merge_3/merged_table.qza \  
430    --i-table2 dada2-denoise_single_1/table.qza \  
    --output-dir feature-table-merge_4  
  
432 # versions: {'feature-table': '2017.4.0'}  
qiime feature-table filter-samples \  
434    --i-table feature-table-merge_4/merged_table.qza \  
    --p-min-frequency 2000 \  
436    --p-min-features 0 \  
    --output-dir feature-table-filter_samples_1  
  
438 # versions: {'feature-table': '2017.4.0'}  
qiime feature-table filter-samples \  
440    --i-table feature-table-filter_samples_1/filtered_table.qza \  
    --m-sample-metadata-file metadata6.txt \  
442    --p-min-frequency 0 \  
    --p-min-features 0 \  
444    --output-dir feature-table-filter_samples_2  
  
# versions: {'taxa': '2017.9.0.dev0+2.g2ebb91d'}  
446 qiime taxa collapse \  
    --i-table feature-table-filter_samples_2/filtered_table.qza \  
448    --i-taxonomy feature-classifier-classify_sklearn_1/classifier.qza \  
    --p-level 6 \  
450    --output-dir taxa-collapse_1  
  
# versions: {'longitudinal': '2018.8.0'}  
452 qiime longitudinal feature-volatility \  
    --i-table taxa-collapse_1/collapsed_table.qza \  
454    --m-metadata-file metadata7.txt \  
    --p-state-column 'month' \  
456    --p-individual-id-column 'studyid' \  
    --p-cv 5 \  
458    --p-n-jobs 4 \  
    --p-n-estimators 100 \
```

```
460 --p-estimator 'RandomForestRegressor' \  
--p-no-parameter-tuning \  
462 --p-missing-samples 'error' \  
--output-dir longitudinal-feature_volatility_1  
464
```

Application programmer interface steps

```
import qiime2  
466 from qiime2.plugins import feature_table  
from qiime2.plugins import feature_classifier  
468 from qiime2.plugins import demux  
from qiime2.plugins import dada2  
470 from qiime2.plugins import longitudinal  
from qiime2.plugins import taxa  
  
472 metadata1_txt_ColumnName =  
qiime2.Metadata.load('metadata1.txt').get_column('ColumnName')  
474 metadata2_txt_ColumnName =  
qiime2.Metadata.load('metadata2.txt').get_column('ColumnName')  
476 metadata3_txt_ColumnName =  
qiime2.Metadata.load('metadata3.txt').get_column('ColumnName')  
478 metadata4_txt_ColumnName =  
qiime2.Metadata.load('metadata4.txt').get_column('ColumnName')  
480 metadata5_txt_ColumnName =  
qiime2.Metadata.load('metadata5.txt').get_column('ColumnName')  
482 metadata6_txt = qiime2.Metadata.load('metadata6.txt')  
metadata7_txt = qiime2.Metadata.load('metadata7.txt')  
  
484 # versions: {'framework': '2017.2.0'}  
99_otu_taxonomy_txt = qiime2.Artifact.import_data(  
486     'FeatureData[Taxonomy]', '99_otu_taxonomy.txt', view_type='TaxonomyFormat'  
)  
  
488 # versions: {'framework': '2017.2.0'}  
99_otus_fasta = qiime2.Artifact.import_data(  
490     'FeatureData[Sequence]', '99_otus.fasta', view_type='DNAFASTAFormat'  
)  
  
492 # versions: {'feature-classifier': '2017.2.0'}  
feature_classifier_fit_classifier_naive_bayes_1 =  
494 feature_classifier.actions.fit_classifier_naive_bayes(  
     reference_reads=99_otus_fasta,  
496     reference_taxonomy=99_otu_taxonomy_txt,  
     classify__alpha=0.01,  
498     classify__chunk_size=-1,  
     classify__class_prior='null',  
500     classify__fit_prior=True,  
     feat_ext__analyzer='char_wb',
```

```
502     feat_ext__binary=False,
      feat_ext__decode_error='strict',
504     feat_ext__encoding='utf-8',
      feat_ext__input='content',
506     feat_ext__lowercase=True,
      feat_ext__n_features=8192,
508     feat_ext__ngram_range=[8, 8]',
      feat_ext__non_negative=True,
510     feat_ext__norm='l2',
      feat_ext__preprocessor='null',
512     feat_ext__stop_words='null',
      feat_ext__strip_accents='null',
514     feat_ext__token_pattern='(?u)\b\w\w+\b',
      feat_ext__tokenizer='null',
516 )

# versions: {'framework': '2017.4.0'}
518 EMPSingleEndDirFmtimport_dir = qiime2.Artifact.import_data(
      'RawSequences', 'EMPSingleEndDirFmtimport_dir', view_type='EMPSingleEndDirFmt'
520 )

# versions: {'demux': '2017.4.0'}
522 demux_emp_single_1 = demux.actions.emp_single(
      seqs=EMPSingleEndDirFmtimport_dir,
524     barcodes=metadata1_txt_ColumnName,
      rev_comp_barcodes=False,
526     rev_comp_mapping_barcodes=True,
      )

528 # versions: {'dada2': '2017.4.0'}
dada2_denoise_single_1 = dada2.actions.denoise_single(
530     demultiplexed_seqs=demux_emp_single_1.per_sample_sequences,
      trunc_len=150,
532     trim_left=0,
      max_ee=2.0,
534     trunc_q=2,
      chimera_method='pooled',
536     min_fold_parent_over_abundance=1.0,
      n_threads=1,
538     n_reads_learn=1000000,
      hashed_feature_ids=True,
540 )

# versions: {'framework': '2017.4.0'}
542 EMPSingleEndDirFmtimport_dir = qiime2.Artifact.import_data(
      'RawSequences', 'EMPSingleEndDirFmtimport_dir', view_type='EMPSingleEndDirFmt'
544 )
```

```
# versions: {'demux': '2017.4.0'}
546 demux_emp_single_2 = demux.actions.emp_single(
    seqs=EMPSingleEndDirFmtimport_dir,
548 barcodes=metadata2_txt_ColumnName,
    rev_comp_barcodes=False,
550 rev_comp_mapping_barcodes=True,
    )

552 # versions: {'dada2': '2017.4.0'}
dada2_denoise_single_2 = dada2.actions.denoise_single(
554 demultiplexed_seqs=demux_emp_single_2.per_sample_sequences,
    trunc_len=150,
556 trim_left=0,
    max_ee=2.0,
558 trunc_q=2,
    chimera_method='pooled',
560 min_fold_parent_over_abundance=1.0,
    n_threads=1,
562 n_reads_learn=1000000,
    hashed_feature_ids=True,
564 )

# versions: {'framework': '2017.4.0'}
566 EMPSingleEndDirFmtimport_dir = qiime2.Artifact.import_data(
    'RawSequences', 'EMPSingleEndDirFmtimport_dir', view_type='EMPSingleEndDirFmt'
568 )

# versions: {'demux': '2017.4.0'}
570 demux_emp_single_3 = demux.actions.emp_single(
    seqs=EMPSingleEndDirFmtimport_dir,
572 barcodes=metadata3_txt_ColumnName,
    rev_comp_barcodes=False,
574 rev_comp_mapping_barcodes=True,
    )

576 # versions: {'dada2': '2017.4.0'}
dada2_denoise_single_3 = dada2.actions.denoise_single(
578 demultiplexed_seqs=demux_emp_single_3.per_sample_sequences,
    trunc_len=150,
580 trim_left=0,
    max_ee=2.0,
582 trunc_q=2,
    chimera_method='pooled',
584 min_fold_parent_over_abundance=1.0,
    n_threads=1,
586 n_reads_learn=1000000,
    hashed_feature_ids=True,
```



```
588 )

# versions: {'framework': '2017.4.0'}
590 EMPSingleEndDirFmtimport_dir = qiime2.Artifact.import_data(
    'RawSequences', 'EMPSingleEndDirFmtimport_dir', view_type='EMPSingleEndDirFmt'
592 )

# versions: {'demux': '2017.4.0'}
594 demux_emp_single_4 = demux.actions.emp_single(
    seqs=EMPSingleEndDirFmtimport_dir,
596 barcodes=metadata4_txt_ColumnName,
    rev_comp_barcodes=False,
598 rev_comp_mapping_barcodes=True,
    )

# versions: {'dada2': '2017.4.0'}
600 dada2_denoise_single_4 = dada2.actions.denoise_single(
602 demultiplexed_seqs=demux_emp_single_4.per_sample_sequences,
    trunc_len=150,
604 trim_left=0,
    max_ee=2.0,
606 trunc_q=2,
    chimera_method='pooled',
608 min_fold_parent_over_abundance=1.0,
    n_threads=1,
610 n_reads_learn=1000000,
    hashed_feature_ids=True,
612 )

# versions: {'framework': '2017.4.0'}
614 EMPSingleEndDirFmtimport_dir = qiime2.Artifact.import_data(
    'RawSequences', 'EMPSingleEndDirFmtimport_dir', view_type='EMPSingleEndDirFmt'
616 )

# versions: {'demux': '2017.4.0'}
618 demux_emp_single_5 = demux.actions.emp_single(
    seqs=EMPSingleEndDirFmtimport_dir,
620 barcodes=metadata5_txt_ColumnName,
    rev_comp_barcodes=False,
622 rev_comp_mapping_barcodes=True,
    )

# versions: {'dada2': '2017.4.0'}
624 dada2_denoise_single_5 = dada2.actions.denoise_single(
626 demultiplexed_seqs=demux_emp_single_5.per_sample_sequences,
    trunc_len=150,
628 trim_left=0,
    max_ee=2.0,
```

```
630     trunc_q=2,  
        chimera_method='pooled',  
632     min_fold_parent_over_abundance=1.0,  
        n_threads=1,  
634     n_reads_learn=1000000,  
        hashed_feature_ids=True,  
636 )  
  
# versions: {'feature-table': '2017.4.0'}  
638 feature_table_merge_seq_data_1 = feature_table.actions.merge_seq_data(  
    data1=dada2_denoise_single_5.representative_sequences,  
640    data2=dada2_denoise_single_4.representative_sequences,  
    )  
  
# versions: {'feature-table': '2017.4.0'}  
642 feature_table_merge_seq_data_2 = feature_table.actions.merge_seq_data(  
    data1=feature_table_merge_seq_data_1.merged_data,  
644    data2=dada2_denoise_single_3.representative_sequences,  
646 )  
  
# versions: {'feature-table': '2017.4.0'}  
648 feature_table_merge_seq_data_3 = feature_table.actions.merge_seq_data(  
    data1=feature_table_merge_seq_data_2.merged_data,  
650    data2=dada2_denoise_single_2.representative_sequences,  
    )  
  
# versions: {'feature-table': '2017.4.0'}  
652 feature_table_merge_seq_data_4 = feature_table.actions.merge_seq_data(  
    data1=feature_table_merge_seq_data_3.merged_data,  
654    data2=dada2_denoise_single_1.representative_sequences,  
656 )  
  
# versions: {'feature-classifier': '2017.5.0'}  
658 feature_classifier_classify_sklearn_1 =  
feature_classifier.actions.classify_sklearn(  
660     reads=feature_table_merge_seq_data_4.merged_data,  
     classifier=feature_classifier_fit_classifier_naive_bayes_1.classifier,  
662     chunk_size=262144,  
     n_jobs=4,  
664     pre_dispatch='2*n_jobs',  
     confidence=0.7,  
666     read_orientation=None,  
    )  
  
# versions: {'feature-table': '2017.4.0'}  
668 feature_table_merge_1 = feature_table.actions.merge(  
670     table1=dada2_denoise_single_5.table,  
     table2=dada2_denoise_single_4.table,
```

```
672 )

# versions: {'feature-table': '2017.4.0'}
674 feature_table_merge_2 = feature_table.actions.merge(
    table1=feature_table_merge_1.merged_table,
676     table2=dada2_denoise_single_3.table,
    )

# versions: {'feature-table': '2017.4.0'}
678 feature_table_merge_3 = feature_table.actions.merge(
    table1=feature_table_merge_2.merged_table,
680     table2=dada2_denoise_single_2.table,
    )
682 )

# versions: {'feature-table': '2017.4.0'}
684 feature_table_merge_4 = feature_table.actions.merge(
    table1=feature_table_merge_3.merged_table,
686     table2=dada2_denoise_single_1.table,
    )

# versions: {'feature-table': '2017.4.0'}
688 feature_table_filter_samples_1 = feature_table.actions.filter_samples(
690     table=feature_table_merge_4.merged_table,
    min_frequency=2000,
692     max_frequency=None,
    min_features=0,
694     max_features=None,
    sample_metadata=None,
696     where=None,
    )

# versions: {'feature-table': '2017.4.0'}
698 feature_table_filter_samples_2 = feature_table.actions.filter_samples(
700     table=feature_table_filter_samples_1.filtered_table,
    sample_metadata=metadata6_txt,
702     min_frequency=0,
    max_frequency=None,
704     min_features=0,
    max_features=None,
706     where=None,
    )

# versions: {'taxa': '2017.9.0.dev0+2.g2ebb91d'}
708 taxa_collapse_1 = taxa.actions.collapse(
710     table=feature_table_filter_samples_2.filtered_table,
    taxonomy=feature_classifier_classify_sklearn_1.classifier,
712     level=6,
    )
```

```
714 # versions: {'longitudinal': '2018.8.0'}
longitudinal_feature_volatility_1 = longitudinal.actions.feature_volatility(
716     table=taxa_collapse_1.collapsed_table,
     metadata=metadata7_txt,
718     state_column='month',
     individual_id_column='studyid',
720     cv=5,
     random_state=None,
722     n_jobs=4,
     n_estimators=100,
724     estimator='RandomForestRegressor',
     parameter_tuning=False,
726     missing_samples='error',
)
728
```

Figure 1c (c-taxa-barplot.qzv in Supplementary File 1)

Data were imported into QIIME 2 as multiplexed 2x150 MiSeq reads and demultiplexed. DADA2¹¹ was applied to single-end reads (as approximately 30% of reads failed to join due to the relatively short sequence length) with no trimming of reads. Taxonomy was assigned to the resulting amplicon sequence variants (ASVs) against the SILVA version 132 99% OTUs (trimmed to the 515F/806R region of the 16S) using q2-feature-classifier's `classify-sklearn` method¹³.

734

Command line interface steps

```
# versions: {'types': '2018.6.0'}
736 qiime tools import \
     --type 'FeatureData[Taxonomy]' \
738     --input-path 7_level_taxonomy.txt \
     --input-format HeaderlessTSVTaxonomyFormat \
740     --output-path 7_level_taxonomy.txt.qza

# versions: {'types': '2018.6.0'}
742 qiime tools import \
     --type 'FeatureData[Sequence]' \
744     --input-path silva132_99.fna \
     --input-format DNAFASTAFormat \
746     --output-path silva132_99.fna.qza

# versions: {'feature-classifier': '2018.6.0', 'types': '2018.6.0'}
748 qiime feature-classifier extract-reads \
     --i-sequences silva132_99.fna.qza \
750     --p-f-primer 'GTGCCAGCMGCCGCGGTAA' \
     --p-r-primer 'GGACTACHVGGGTWTCTAAT' \
752     --p-trunc-len 0 \
     --p-trim-left 0 \
754     --p-identity 0.8 \
```

```
--output-dir feature-classifier-extract_reads_1

756 # versions: {'feature-classifier': '2018.6.0', 'types': '2018.6.0'}
qiime feature-classifier fit-classifier-naive-bayes \
758 --i-reference-reads feature-classifier-extract_reads_1/reads.qza \
--i-reference-taxonomy 7_level_taxonomy.txt.qza \
760 --p-classify--alpha 0.001 \
--p-classify--chunk-size 20000 \
762 --p-classify--class-prior 'null' \
--p-no-classify--fit-prior \
764 --p-no-feat-ext--alternate-sign \
--p-feat-ext--analyzer 'char_wb' \
766 --p-no-feat-ext--binary \
--p-feat-ext--decode-error 'strict' \
768 --p-feat-ext--encoding 'utf-8' \
--p-feat-ext--input 'content' \
770 --p-feat-ext--lowercase \
--p-feat-ext--n-features 8192 \
772 --p-feat-ext--ngram-range '[7, 7]' \
--p-no-feat-ext--non-negative \
774 --p-feat-ext--norm 'l2' \
--p-feat-ext--preprocessor 'null' \
776 --p-feat-ext--stop-words 'null' \
--p-feat-ext--strip-accents 'null' \
778 --p-feat-ext--token-pattern '(?u)\b\w\w+\b' \
--p-feat-ext--tokenizer 'null' \
780 --output-dir feature-classifier-fit_classifier_naive_bayes_1

# versions: {'types': '2018.8.0'}
782 qiime tools import \
--type 'SampleData[SequencesWithQuality]' \
784 --input-path se-64-manifest.csv \
--input-format SingleEndFastqManifestPhred64 \
786 --output-path se-64-manifest.csv.qza

# versions: {'quality-filter': '2018.8.0', 'types': '2018.8.0'}
788 qiime quality-filter q-score \
--i-demux se-64-manifest.csv.qza \
790 --p-min-quality 4 \
--p-quality-window 3 \
792 --p-min-length-fraction 0.75 \
--p-max-ambiguous 0 \
794 --output-dir quality-filter-q_score_1

# versions: {'types': '2018.8.0', 'deblur': '2018.8.0'}
796 qiime deblur denoise-16S \
--i-demultiplexed-seqs quality-filter-q_score_1/filtered_sequences.qza \
798 --p-trim-length 85 \
```

```
800 --p-no-sample-stats \  
--p-mean-error 0.005 \  
802 --p-indel-prob 0.01 \  
--p-indel-max 3 \  
804 --p-min-reads 10 \  
--p-min-size 2 \  
--p-jobs-to-start 8 \  
806 --p-hashed-feature-ids \  
--output-dir deblur-denoise_16S_1  
  
808 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
qiime feature-table filter-samples \  
810 --i-table deblur-denoise_16S_1/table.qza \  
--m-metadata-file metadata1.txt \  
812 --p-min-frequency 2000 \  
--p-min-features 0 \  
814 --p-where "Site='Steep'" \  
--p-no-exclude-ids \  
816 --output-dir feature-table-filter_samples_1  
  
# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
818 qiime feature-table filter-features \  
--i-table feature-table-filter_samples_1/filtered_table.qza \  
820 --p-min-frequency 100 \  
--p-min-samples 0 \  
822 --p-no-exclude-ids \  
--output-dir feature-table-filter_features_1  
  
824 # versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}  
qiime feature-table filter-seqs \  
826 --i-data deblur-denoise_16S_1/representative_sequences.qza \  
--i-table feature-table-filter_features_1/filtered_table.qza \  
828 --p-no-exclude-ids \  
--output-dir feature-table-filter_seqs_1  
  
830 # versions: {'feature-classifier': '2018.8.0', 'types': '2018.8.0'}  
qiime feature-classifier classify-sklearn \  
832 --i-reads feature-table-filter_seqs_1/filtered_data.qza \  
--i-classifier feature-classifier-fit_classifier_naive_bayes_1/classifier.qza \  
834 --p-reads-per-batch 0 \  
--p-n-jobs 8 \  
836 --p-pre-dispatch '2*n_jobs' \  
--p-confidence 0.7 \  
838 --output-dir feature-classifier-classify_sklearn_1  
  
# versions: {'types': '2018.11.0.dev0', 'feature-table':  
840 '2018.11.0.dev0+3.g345ac6e'}  
qiime feature-table group \  

```

```
842 --i-table feature-table-filter_features_1/filtered_table.qza \  
--m-metadata-file metadata2.txt \  
844 --m-metadata-column 'ColumnName' \  
--p-axis 'sample' \  
846 --p-mode 'median-ceiling' \  
--output-dir feature-table-group_1  
  
848 # versions: {'taxa': '2018.11.0.dev0', 'types': '2018.11.0.dev0'}  
qiime taxa barplot \  
850 --i-table feature-table-group_1/grouped_table.qza \  
--i-taxonomy feature-classifier-classify_sklearn_1/classification.qza \  
852 --m-metadata-file metadata3.txt \  
--output-dir taxa-barplot_1  
854
```

Application programmer interface steps

```
import qiime2  
856 from qiime2.plugins import feature_table  
from qiime2.plugins import feature_classifier  
858 from qiime2.plugins import taxa  
from qiime2.plugins import quality_filter  
860 from qiime2.plugins import deblur  
  
metadata1_txt = qiime2.Metadata.load('metadata1.txt')  
862 metadata2_txt_ColumnName =  
qiime2.Metadata.load('metadata2.txt').get_column('ColumnName')  
864 metadata3_txt = qiime2.Metadata.load('metadata3.txt')  
  
# versions: {'types': '2018.6.0'}  
866 7_level_taxonomy_txt = qiime2.Artifact.import_data(  
    'FeatureData[Taxonomy]', '7_level_taxonomy.txt',  
868 view_type='HeaderlessTSVTaxonomyFormat'  
)  
  
870 # versions: {'types': '2018.6.0'}  
silva132_99_fna = qiime2.Artifact.import_data(  
872 'FeatureData[Sequence]', 'silva132_99.fna', view_type='DNAFASTAFormat'  
)  
  
874 # versions: {'feature-classifier': '2018.6.0', 'types': '2018.6.0'}  
feature_classifier_extract_reads_1 = feature_classifier.actions.extract_reads(  
876 sequences=silva132_99_fna,  
f_primer='GTGCCAGCMGCCGCGGTAA',  
878 r_primer='GGACTACHVGGGTWTCTAAT',  
trunc_len=0,  
880 trim_left=0,  
identity=0.8,  
882 )
```

```
# versions: {'feature-classifier': '2018.6.0', 'types': '2018.6.0'}
884 feature_classifier_fit_classifier_naive_bayes_1 =
feature_classifier.actions.fit_classifier_naive_bayes(
886     reference_reads=feature_classifier_extract_reads_1.reads,
reference_taxonomy=7_level_taxonomy_txt,
888     classify__alpha=0.001,
classify__chunk_size=20000,
890     classify__class_prior='null',
classify__fit_prior=False,
892     feat_ext__alternate_sign=False,
feat_ext__analyzer='char_wb',
894     feat_ext__binary=False,
feat_ext__decode_error='strict',
896     feat_ext__encoding='utf-8',
feat_ext__input='content',
898     feat_ext__lowercase=True,
feat_ext__n_features=8192,
900     feat_ext__ngram_range='[7, 7]',
feat_ext__non_negative=False,
902     feat_ext__norm='l2',
feat_ext__preprocessor='null',
904     feat_ext__stop_words='null',
feat_ext__strip_accents='null',
906     feat_ext__token_pattern='(?u)\b\w\w+\b',
feat_ext__tokenizer='null',
908 )

# versions: {'types': '2018.8.0'}
910 se-64-manifest_csv = qiime2.Artifact.import_data(
'SampleData[SequencesWithQuality]', 'se-64-manifest.csv',
912 view_type='SingleEndFastqManifestPhred64'
)

# versions: {'quality-filter': '2018.8.0', 'types': '2018.8.0'}
914 quality_filter_q_score_1 = quality_filter.actions.q_score(
quality_filter_q_score_1 = quality_filter.actions.q_score(
916     demux=se-64-manifest_csv,
min_quality=4,
918     quality_window=3,
min_length_fraction=0.75,
920     max_ambiguous=0,
)

# versions: {'types': '2018.8.0', 'deblur': '2018.8.0'}
922 deblur_denoise_16S_1 = deblur.actions.denoise_16S(
deblur_denoise_16S_1 = deblur.actions.denoise_16S(
924     demultiplexed_seqs=quality_filter_q_score_1.filtered_sequences,
trim_length=85,
926     sample_stats=False,
```



```
mean_error=0.005,
928 indel_prob=0.01,
    indel_max=3,
930 min_reads=10,
    min_size=2,
932 jobs_to_start=8,
    hashed_feature_ids=True,
934 )

# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}
936 feature_table_filter_samples_1 = feature_table.actions.filter_samples(
    table=deblur_denoise_16S_1.table,
938 metadata=metadatal_txt,
    min_frequency=2000,
940 max_frequency=None,
    min_features=0,
942 max_features=None,
    where="Site='Steep'",
944 exclude_ids=False,
)

# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}
946 feature_table_filter_features_1 = feature_table.actions.filter_features(
948 table=feature_table_filter_samples_1.filtered_table,
    min_frequency=100,
950 max_frequency=None,
    min_samples=0,
952 max_samples=None,
    metadata=None,
954 where=None,
    exclude_ids=False,
956 )

# versions: {'types': '2018.8.0', 'feature-table': '2018.8.0'}
958 feature_table_filter_seqs_1 = feature_table.actions.filter_seqs(
    data=deblur_denoise_16S_1.representative_sequences,
960 table=feature_table_filter_features_1.filtered_table,
    metadata=None,
962 where=None,
    exclude_ids=False,
964 )

# versions: {'feature-classifier': '2018.8.0', 'types': '2018.8.0'}
966 feature_classifier_classify_sklearn_1 =
feature_classifier.actions.classify_sklearn(
968 reads=feature_table_filter_seqs_1.filtered_data,
    classifier=feature_classifier_fit_classifier_naive_bayes_1.classifier,
970 reads_per_batch=0,
```

```
972     n_jobs=8,  
     pre_dispatch='2*n_jobs',  
     confidence=0.7,  
974     read_orientation=None,  
 )  
  
976 # versions: {'types': '2018.11.0.dev0', 'feature-table':  
 '2018.11.0.dev0+3.g345ac6e'}  
978 feature_table_group_1 = feature_table.actions.group(  
     table=feature_table_filter_features_1.filtered_table,  
980     metadata=metadata2_txt_ColumnName,  
     axis='sample',  
982     mode='median-ceiling',  
 )  
  
984 # versions: {'taxa': '2018.11.0.dev0', 'types': '2018.11.0.dev0'}  
taxa_barplot_1 = taxa.actions.barplot(  
986     table=feature_table_group_1.grouped_table,  
     taxonomy=feature_classifier_classify_sklearn_1.classification,  
988     metadata=metadata3_txt,  
 )  
990
```

Figure 1d (d-ili-plot.qzv in Supplementary File 1)

992 The input files for this visualization are a stereolithography file (STL) and a sample metadata file with a
mapping between samples and the spatial coordinates (x, y and z). Both files were obtained from `ili's [GitHub](#)
994 page^{15,16}. The comma-separated file was converted into a tab-separated format (to make it compatible with
QIIME 2).

Command line interface steps

```
996 # versions: {'ili': 'v0.1.1'}  
qiime tools import \  
998     --type 'Model' \  
     --input-path model.stl \  
1000     --input-format STLFile \  
     --output-path model.stl.qza  
  
1002 # versions: {'ili': 'v0.1.1'}  
qiime ili plot \  
1004     --i-model model.stl.qza \  
     --m-metadata-file metadata1.txt \  
1006     --output-dir ili-plot_1
```

Application programmer interface steps

```
1008 import qiime2  
from qiime2.plugins import ili
```

```
010  metadata1_txt = qiime2.Metadata.load('metadata1.txt')
    # versions: {'ili': 'v0.1.1'}
012  model_stl = qiime2.Artifact.import_data(
    'Model', 'model.stl', view_type='STLFile'
014  )
    # versions: {'ili': 'v0.1.1'}
016  ili_plot_1 = ili.actions.plot(
    model=model_stl,
018  metadata=metadata1_txt,
    )
```

Supplementary File 3: Notes from end users on QIIME 2

We solicited feedback from a broad contingent of QIIME 2 users and developers (N = 15), asking them to review different aspects of QIIME 2 that they use in their research. The majority of these reviewers are independent users whom we selected on the basis of their activity on the QIIME 2 Forum. Four reviewers (identified below) are co-authors of this article who contributed to QIIME 2 (e.g., in the form of plugins compatible with QIIME 2, documentation, or other major contributions), but are currently unaffiliated with the Caporaso or Knight research groups (the two groups who derive direct funding from the primary QIIME 2 NSF grant). Their use of QIIME 2 is what initiated involvement in QIIME 2 development, and hence we feel that these investigators can provide unique insight on how QIIME 2 attracts a diverse community of users and developers. Reviewers are listed in alphabetical order. All reviewers gave permission to be quoted in this document, and approved the inclusion of their statements as they appear here (e.g., in some cases, typos were corrected or statements were abridged).

Some end-users have commented specifically on plugins that are included in the worked examples in the text (Fig 1), including q2-emperor (Fig 1a) and q2-longitudinal (Fig 1b).

Ahmed Abdelfattah, PhD.

Postdoctoral Fellow

Department of Ecology, Environment and Plant Sciences, Stockholm University, Sweden

QIIME is by far my favorite pipeline for microbial data analysis. I started using it while I was doing my PhD. And I have not stopped since. While very comprehensive, it is user-friendly, even for those who have little experience in bioinformatics. QIIME is a perfect example of open source software, making it both versatile and able to cope with the ever-expanding field of microbial ecology. QIIME 2 takes microbial data analysis a big step further. While keeping the main essential features of QIIME 1, the newly added plugins are very useful. Perhaps my favorite plugin is q2-longitudinal, a plugin for longitudinal and paired-sample analyses. This plugin contains pipelines, methods, and visualizers, all of which provide complex analysis to evaluate paired samples over a period of time.

Rozlyn C.T. Boutin, BScH

MD/PhD student

University of British Columbia, Vancouver, Canada

I have been working with a range of features available in QIIME 2, including many of those described in the "Moving Pictures Tutorial". I have found each of these plugins straightforward to use and well described in the online tutorial. I found it especially useful that the tutorial and QIIME 2 websites describe what each plugin does, what parameters in each analysis can be changed and how, as well as provide guidance regarding how to modify the parameters of a command to more appropriately fit the study in question.

More recently, I have been using the q2-longitudinal plugin to compute a gut microbiome "maturity index" for stool samples collected at multiple time points from the same subjects and ultimately compare how the maturity of the gut microbiome differs between subjects with different disease states... The ability to perform analyses comparing gut microbiota maturity between groups of individuals with differing disease susceptibilities demonstrates the unique ability of QIIME 2 to perform high-level bioinformatic and statistical analyses while also taking complex biological contexts and ecological factors into consideration. Indeed,

40 QIIME 2 does a great job of bridging the gaps in language and scientific expertise between the disciplines of
42 biology, bioinformatics, and statistics. Moreover, the available plugins accommodate novel, complex
analyses being implemented in the latest literature.

44 While R code for gut microbiome maturity index predictions is available, the QIIME 2 plugin provides an
easy-to-use platform with customizable input and output options. User-friendly tuning parameters, such as
46 the ability to add spaghetti vectors to volatility charts, modify the fraction of samples used for training vs.
testing of the regressor, stratify training and test data among metadata categories, and set a seed to ensure
48 repeatability among users working with the same data are all useful and relevant features of the plugin that
are easily and intuitively manipulated... Although it is not possible to obtain a microbiota-by-age Z (MAZ)
50 score for samples used to train the model, the output of the q2-longitudinal maturity index prediction
provides both intuitive visualization output files as well as the option of downloading a TSV file containing a
52 MAZ score for samples not used for training, which can subsequently be used in downstream analyses.
Allowing for yet more advanced downstream analyses, the plugin also provides outputs with detailed
54 information on how to interpret the output of the command; for instance, the feature-importance scores can
be exported to determine which features are most important for discriminating microbial maturity indices, and
56 how each feature changes over time in each experimental group can be visualized in the cluster map output.
Each of these output files are intuitive to interpret and the resulting visualization files aesthetically pleasing.

58 All of the plugins I have used in QIIME 2 thus far have been easy to use and I have been impressed by the
advanced and innovative analyses available. These tools are complemented with a well maintained user
60 website and easy-to-follow tutorials that provide helpful guidance without burdening the reader with too
many details. Details, however, are still available elsewhere on the website for those who are interested and
the QIIME 2 forum is exceedingly helpful and responsive.

62 **David J. Bradshaw II**

Ph.D. Student

64 **Florida Atlantic University - Harbor Branch Oceanographic Institute**

66 QIIME 2 has been essential to my PhD work, allowing me an easy to use and very flexible way to analyze
my 16S rRNA amplicon sequences. My colleagues and I also plan on using this system for 18S analysis. I
68 have been a user of QIIME 2 for about a year now, and I have been able to teach myself how to use it with
the assistance of their well written guides, along with some very timely and useful advice from creators and
users. I do not have access to a microbiome expert at my school, so these resources have been very
70 important in helping me complete my work. QIIME 2's system of plugins basically allows me to do anything
that I need to do in order to make sure that I have the best quality sequences for analysis and also provides
72 many ways to analyze the data for my dissertation and publications. There are so many different ways to
analyze microbiome data, and QIIME 2 is robust enough to allow you to analyze it however you prefer, such
74 as using the Deblur, vsearch, or DADA2. **[note from the QIIME 2 authors: these are QIIME 2 plugins
available to perform various actions for denoising, dereplication, and/or OTU clustering.]**

76 **Lorinda Bullington**

Molecular Ecologist

78 **MPG Ranch, Missoula, Montana**

80 The developers of the QIIME 2 bioinformatics platform foster an interactive, collaborative environment
wherein users can directly influence the functional ability of the platform. The versatility of the many plugins
allows researchers to use a single platform to analyze amplicon data from various primer pairs and multiple

82 target organisms including bacteria, fungal endophytes, and arbuscular mycorrhizal fungi. The QIIME 2
84 forum provides quick and thorough answers to questions and allows direct communication between users
and developers which greatly increases QIIME 2's accessibility for early career scientists. Data provenance
86 tracking ensures reproducible results and is very helpful when returning to old datasets or merging new
ones. These aspects, along with the excellent data visualizations associated with each step help to eliminate
the black box that is so often associated with bioinformatics pipelines.

88 **Justine W. Debelius, Ph.D.**

Postdoctoral Scholar

90 **Karolinska Institutet Dept. of Medical Epidemiology and Biostatistics, Stockholm, Sweden**

Dr. Debelius is not an author of this paper, but worked previously in the laboratory of Dr. Rob Knight and has taught
92 QIIME 1 workshops.

As someone who considers herself somewhere on the advanced to expert level of user, I love the capacity
94 [in QIIME 2] to have multiple interfaces. I appreciate that I can generate an object in QIIME 2 and then use a
standard python API to pull it into a notebook for my own modifications or perform more complex operations.
96 It means that I can easily perform more complex analyses (like specialised analyses with Gneiss) that aren't
possible via the command line interface.

I'm also a huge fan of QIIME 2 view: it makes it super easy to share results with my PI and collaborators.
98 We're discussing including a QIIME 2 view object in the supplement of a paper I'm a co-author of, because it
100 makes the interaction easier.

Finally, I like the more "choose your own adventure" aspect of the plugin architecture. Rather than trying to
102 structure a parameters file in QIIME 2, it's much easier to choose the best or preferred algorithm for each
step. I can even mix tools: paired end joining in vsearch, denoising in Deblur, and chimera slaying with
104 another algorithm...

Claire Duvall

106 **Ph.D. Candidate**

Dept. of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA

108 Dr. Duvall is a co-author of this paper.

My favorite (non-obvious) thing about QIIME 2 is that it provides a really easy "in" for people to join the ranks
110 of the "legit" microbiome research community. I think this works in two big ways:

- 112 1. Researchers who work at institutions without many resources now have access to the latest and
greatest tools in our field. I love that QIIME 2 does a lot to even the playing field and open up this
114 field to many more researchers. There are many axes to this: scientists in developing countries,
clinicians without much data background, citizen scientists who got their microbiomes sequenced,
etc — all now have access where before they didn't! QIIME 2 is unique in this way over other
116 bioinformatics software because there are so many different ways to engage with it, and it was
obviously created with this impact in mind.
- 118 2. Established researchers now get to add their "school of thought" to the mainstream conversation...
Now, with QIIME 2, I get to write in my school of thought. I can write plugins for my favorite methods,
120 add parameters that are missing in existing wrappers, give my perspective on how things "should" be
done on the forum, and incorporate my school of thought into sort-of canon via writing tutorials or
122 documentation... Now, if you don't like what the field is doing, you can actively change it!

A few other more specific experiences from my own interactions with QIIME 2:

1. The plugin philosophy makes it suuuuper easy to get any new method I make into as many hands as possible, and really incentivizes making code more useable. For example, I helped a postdoc in our lab develop a new method to normalize data to reduce batch effects, and was really excited when Greg [Caporaso] reached out and asked if we wanted to add our method into QIIME 2. The common framework that many people use and which provided a clear approach for how to convert our scripts into a bona-fide tool was the nudge I needed to convince me to put in the extra effort to turn our scripts into something more useable (the method is now available in QIIME 2 as q2-perc-norm). Everybody wins here: the community has a new method that is easy to install and use, and I get to share my method broadly and also say that I've contributed to open-source software development.
2. Although I don't personally use QIIME 2 that much (I'm past most of the heavy data-processing parts of my PhD), I tell all my peers to use it if only because of the convenience of the platform: having one standard framework that wraps existing functions across varying software packages and tools is so nice because installation issues aren't such huge obstacles and you don't have to constantly write intermediate scripts to reformat input and output files. It's just so nice to have one bash script that can go from raw data to feature table in a really streamlined way.
3. Before I "converted" to QIIME 2, I was the main person in charge of maintaining and developing our lab's in-house data processing pipeline. This took me between a few hours to a few days per month, ranging from helping lab members learn to use the pipeline to fixing bugs to implementing new functionalities. Now that QIIME 2 exists, I can point lab members to the forum for help, raise issues for the core developers to address when I find difficult bugs (or solve them myself if they are easy to address), and spend the rest of that time saved on either my own research or on adding new functionality to QIIME 2 that is unique from my lab. As a specific example, I'm pretty sure it took me more time to incorporate distribution-based OTU clustering into my lab's pipeline (which serves a maximum of 20 people) than it did to write it up as a plugin (which can theoretically serve thousands of people as q2-dbotu). It's clear that joining the open-source QIIME 2 community is a much better investment of everybody's time!

Erika Korzune Ganda, DVM, Ph.D.

Postdoctoral Associate

Department of Food Science, Cornell University, Ithaca, NY

I could not have finished my PhD without the software and the amazing help of the developers through the forum...

I love how I can save my code and be able to go back to it six months later and not only know exactly what I did, but also be able to analyze a new dataset with much more ease. And it is at this point that QIIME 2 offers an amazing advantage with provenance tracking and UUID identifiers for reproducible identification of sequence variants.

In addition to making analysis more straightforward and traceable, the ability to share visualization files over email with collaborator is an AMAZING feature. I am a veterinarian that knows a little about coding and microbiome research. Now I am able to collaborate with clinicians working on data on the microbiome of bovine rumen, companion animals' mouths with periodontal disease, milk, salmon and trout skin, and murine feces. And the best of it is that I can email a .qzv [QIIME 2 visualization] file with some explanation on variable names, and they can then look at the interactive interface without having to go through any coding (which can be pretty scary, if you are not used to it). The web-based interface makes our meetings and collaborations much more effective.

Another great thing is that we are not hostage to the outputs of view.qiime2.org. I can, at any point of my analysis, extract my raw data and make plots using my software of preference.

However, it is not enough to have a good interface without good people working behind the scenes. I give kudos to the developer team that answered my endless questions and helped me through numerous error messages from 2016 to today...

In terms of data analysis, it is no secret that microbiome research is a rapidly evolving field. It is important to keep up to date with the most appropriate analysis methods, as there is still a lot of debate in terms of what the most appropriate data handling and analysis techniques are. This is another advantage of the plugin-based architecture of QIIME 2. It allows for people working on the forefront of data analysis to make their tools available to people like me, who wouldn't necessarily come across them if they weren't available in a user-friendly platform.

Dr. Alexander Mahnert

Postdoctoral researcher

Department of Internal Medicine, Medical University of Graz, Austria

I'm now working with QIIME 2 since the Virtual box image of version 2.0.6. Back then I was already an experienced user of QIIME 1 due to regular analysis of amplicon and shotgun data. I got very excited about QIIME 2 when I realized its main improvements over QIIME 1.9.1: 1st, denoising instead of clustering, 2nd, defined semantic types for different data formats, and 3rd, classifiers to assign taxonomy. We already used denoising algorithms like DADA2 before they were implemented in QIIME since we realized that they usually give us a more conservative view on the profile (e.g., OTUs, RSVs, ASVs) of human-associated Archaea¹. This was a first argument to use QIIME 2 for our data analysis so that we could circumvent inconvenient switching from QIIME 1.9.1 to R and back again. Another argument was my curiosity about the clear core concept of QIIME 2 with artifacts, defined semantics, plugins, methods and visualizers. This clear concept not only guarantees the correct usage of a certain plugin or method, but also helps to keep track of a bunch of qza and qzv files during your analysis by the possibility to check the provenance of each file. In addition, it also saves a lot of hard disk space in the frame of an analysis. And finally, our study² on the microbial dynamics in an isolated and confined built environment (ICE) and the chance to apply plugins like q2-longitudinal on our >500 day observation of some built environment surfaces sealed my transition to QIIME 2 (although we had to publish the manuscript without the longitudinal analysis of QIIME 2 due to time limitations). I'm a big fan of this plugin and its methods like pairwise difference and distance, linear mixed effect models, first rate of change and especially the volatility analysis. All is so well documented and provides a lot of possibilities to visualize your data. And after several improvements to the error messages of certain methods, I think even beginners can make a good analysis of their longitudinal datasets. Now, after the use of QIIME 2 for almost 2 years now, I think it is a good moment to make a first conclusion about the pro's and con's:

First of all, QIIME 2 (as QIIME 1) has a great documentation with elaborate tutorials to give users a step by step introduction to certain plugins (for beginners – <https://docs.qiime2.org/2018.11/tutorials/overview/> or experienced users – <https://docs.qiime2.org/2018.11/tutorials/qiime2-for-experienced-microbiome-researchers/>). With example data, many figures and tables, the developers do not only make sure that users apply their tools for the right purpose, but also help them to understand a tool and to learn how to interpret their own results. And if something is still unclear or buggy, the active QIIME 2 forum is a good address to solve many obstacles.

210 However, as I now realize that my lines about QIIME 2 are pretty positive I also want to indicate several
211 places for improvements. Not all plugins are covered by such great tutorials as in q2-longitudinal, q2-gneiss,
212 q2-quality-control, q2-feature-classifier and q2-sample-classifier. It would be great to see more elaborate
213 tutorials also for other plugins [e.g., q2-composition and q2-picrust2]. I'm a big fan of interactive plots like
214 q2-emperor and [the volatility action in q2-longitudinal]. I think this interactivity is a big part of the fun you
215 experience while using QIIME 2 for your data analysis. Therefore, more options to adjust simple things like
216 colors or export each generated plot also as a svg file would be a nice add in my eyes. Depending on your
217 input data, it might be still necessary to use QIIME 1 before data import to QIIME 2. I think this is a nasty
218 peculiarity. Why isn't there a plugin in QIIME 2 yet to extract barcodes from fastq files to allow direct
compatibility with the available import options? ...

220 To come to an end – I think QIIME 2 is one of the best free packages to analyze amplicon data at the
221 moment and already quite complete. For sure there are still things to update and improve and always will be,
222 but from my point of view there is no sense to teach students still how they could work with QIIME 1.9.1.
223 when there is QIIME 2 available. There are already a few publications out there where we used QIIME 2 for
224 the data analysis³ and there are plenty waiting to be wrapped up into a manuscript. I even like to re-analyze
old data processed in QIIME 1 and use QIIME 2 to get a new perspective on it.

226 Finally, I want to thank the developers of QIIME for QIIME 2, that this software is free and not commercial,
and encourage all of them to keep on improving methods and adding tools to the Quantitative Insights Into
Microbial Ecology!

228 **Dr. Melanie C Melendrez**
229 **Computational Microbiologist**
230 **St. Cloud State University, St. Cloud, MN**

232 One of the best aspects of the QIIME 2 workflow is provenance and command consistency. While the
233 learning curve for QIIME could be considered a bit of a climb, once you become familiar with syntax and if
234 you aren't shy about utilizing and searching the forum, the curve is less of a climb and more like a steady
walk uphill that is not unpleasant.

236 I think it's fantastic the developers of QIIME have recognized the need for consistent tools in the area of
237 metabolomics, metatranscriptomics, and metabolic pathway inference. There's such an inundation in the
238 literature recently of microbiome studies and investigators are under increasing pressure to ensure their
239 study doesn't become a J.A.M - (just another microbiome) study/session. The challenge in expanding
240 beyond 'who is there' is the lack of correspondence and formats required by all the different tools that would
241 have to be navigated; and, in the end, each tool may have a different sequence quality control workflow, a
242 different taxonomic characterization database requirement (or workflow) or a different format requirement -
243 forcing the investigator or manipulate the data in a separate framework that may introduce errors or
244 inconsistency. These are examples of differences that would then potentially greatly impact more complex
245 downstream analyses involving metabolic pathway inferences, metabolomics or metatranscriptomics. By
246 developing plugins for the QIIME 2 framework to address these downstream, more complex analyses, all
247 data stays in the same framework, consistent formats that can be exported to known formats or shared
248 directly through QIIME 2 view, is subject to the same quality control protocols (which are adjustable),
analyzed using the same databases and versions of databases. It greatly enhances and streamlines the
ability to do reproducible science. In short, it's neat.

250 I am also looking forward to q2-metaphlan2 and q2-SCNIC as I have used metaphlan in the past and I am
moving toward network analysis in the future. **[note from the QIIME 2 authors: these are plugins for**

shotgun metagenome taxonomic profiling and network analysis, respectively, which will facilitate multi-omics analysis in QIIME 2.]

I am particularly fond of QIIME 2 view. It is very useful as I have worked with a variety of studies and it is great to have the option of sharing my results without requiring those on the receiving end to figure out how to get the program, install it, and run it in order to see the results I generated.

I am really good at breaking things... and the QIIME forum has become an indispensable tool for me. The QIIME platform is highly utilized in part because of its fantastic user support. If I have issues with linux I go to stackoverflow typically. For various bioinformatics tools they often have no user support beyond Github issues or a Google group. They are useful but what is compelling about using QIIME as a scientist is the non-judgemental, seriously-we-want-to-help-you vibe of a forum dedicated to this platform... QIIME developers and the forum community as a whole are committed, timely with their responses to requests, and have been very friendly from the most basic syntax mistakes to the more complex conceptual or statistical error requests. They are also open to assisting in results interpretation to ensure you are using the plugins correctly.

Currently I am educator that builds all my informatics teaching modules from scratch - I would incorporate QIIME 2 in my teaching as a basic platform to learn about metagenomics and teach students how to analyze microbiome data as well as teach the value of reproducibility and consistency in study design and analysis. I also use it as an example of 'how' scientific inquiry incorporating computational biology should be done — collaboratively; we are all works in progress, moving forward in a field that changes rapidly.

Devon O'rourke
Ph.D. Student
University of New Hampshire

There are a lot of technical reasons to love QIIME 2, but most of my reasons have to do with collaboration. Provenance is a game changer when working in groups - it's exactly the data management version history platform you'd dream up (and greatly simplifies citation tracking to boot). There's obviously a diversity of tools to QC data, merge projects, explore suites of diversity metrics and visualizations, but the biggest asset in terms of collaboration is the ease with which an artifact file is shared - just provide access to a file and you can view it on a browser on any machine. My favorite part though is the broader community with which the QIIME ecosystem exists - when you're thinking about what tools you can depend on, reach out for help with, and possibly contribute back to, there's nothing like it. The documentation is superb, but the community forum (and its members) is ultimately what propels me through the challenges I routinely face.

Adam R. Rivers, Ph.D.
Computational Biologist
USDA-ARS, Gainesville, FL

Dr. Rivers is a co-author of this paper.

I recently developed a stand-alone software package, ITSxpress, to trim fungal sequences for amplicon sequencing analysis and decided to extend it to integrate into QIIME 2. Working with the QIIME 2 community has been a very positive experience. QIIME is a very popular bioinformatics package. In the 8 years since its initial release the paper for QIIME 1 has had over 12,700 citations⁴. Creating software for this preexisting user base has helped my software gain much broader adoption. Since introducing the QIIME 2 plugin several months ago I've had 473 downloads of the Bioconda package for ITSxpress. Having a larger user

base has also helped me improve the quality of the software. QIIME 2 users have identified several bugs that I have been able to rapidly fix. QIIME 2 has a very active user forum and I think that has encouraged users to report issues to me that they may not have reported if they had to email me as the maintainer of a stand-alone package.

QIIME 2 has a very professional design that allows for data provenance tracking, semantic typing and a nice API and good file checking functionalities. The software introduces a number of concepts that you need to understand before writing a package. The documentation for developers is improving, but some more explanation about data types and transformers would be helpful. Fortunately, the QIIME 2 development team is extremely helpful and welcoming to new contributors. I intend to develop more tools for QIIME 2 in the future.

Biswarup Sen

Associate Professor

Tianjin University, School of Environmental Science and Engineering, Tianjin, China

With the rapid growth in the application of high-throughput sequence (HTS) data in contemporary research, the need of a multitasking bioinformatics platform becomes inevitable. The contribution of QIIME 2 as one of the most favorite bioinformatics platforms is unmeasurable in the advancement of HTS data analysis. As a QIIME 2 user and a member of the QIIME 2 forum, I have found QIIME 2 as the most user-friendly and versatile platform that can perform a myriad of analyses rapidly and with much less effort. The tutorials that are provided by the QIIME 2 developers are comprehensive and one with a minimal training on executing commands can begin to use this platform on the chosen operating system. I found the QIIME 2 forum very helpful in troubleshooting the issues that are often faced while running the QIIME 2 plugins. The workflow designed by the developers is well-structured and easy to grasp for initiating data analysis starting from the import of raw sequence data files to taxonomy assignment. QIIME 2 is applicable for analysis of both prokaryotic and eukaryotic communities hailing from any environment, known or underexplored. Besides, the opportunity for other developers to add new plugins allows the creation of useful plugins that are task-specific, for example the ITSxpress plugin can do a quick trimming of the flanking regions in the fungal ITS1/ITS2 sequences, which results in accurate taxonomic assignments down to species level. The ability of post-analysis with the QIIME 2 artifacts on other platforms, e.g. R, makes QIIME 2 more user-friendly and allows the user to perform all statistical analysis and generate high-quality and intuitive plots ready for publication. The developers of QIIME 2 are indeed performing a huge task of providing such a versatile platform which undergoes continuous evaluation, troubleshooting, and upgradation in order to provide the user the ease to conduct their analysis without much challenges and hurdles.

It would take pages for me to describe the various advantages and the variety of tools and plugins wrapped in QIIME 2 platform. In conclusion, I would strongly recommend QIIME 2 for the modern data scientists and young researchers working on HTS data analysis. Appended are a couple of my research group's latest articles where the QIIME platform was used^{5,6}. A few others are in the pipeline.

Solveig Tangedal, M.D.

Dept. of Thoracic Medicine, Haukeland University Hospital, Bergen, Norway

Ph.D. Candidate, Dept. of Clinical Science, Faculty of Medicine and Dentistry, University of Bergen, Norway

I am a medical doctor, and a PhD candidate studying microbiota in airways disease. QIIME 1 was our preferred pipeline for the first article published on microbiota⁷, and it was natural to use QIIME 2 for our next paper.

336 With the transition from QIIME 1 to QIIME 2, several useful tools have been made available. The integration
of DADA2 as a wrapper including also chimera removal has provided a more accurate separation of
338 sequences into units for analyses (amplicon sequence variants) compared to QIIME 1 operational taxonomic
units. Different visualization tools like the quality plots improve the understanding of how sequencing data
340 are processed before accepting sequences for final analyses. The pipeline also allows for better control of
how processes change the datasets during different filtering steps. New statistical, analytical tools are
342 implemented for longitudinal data and for evaluating dynamics in taxonomic compositions in samples. For
my data this could be further improved with an even more comprehensive development for paired analyses
344 with only two timepoints. This would also improve analyses for many other medical research studies with two
sampling points.

346 The output data can be made available for other analytical tools like R... This ensures that even if certain
tools are not made available directly in QIIME 2, the data derived from the pipeline still can be analyzed with
appropriate statistical tests.

348 For us it is especially useful to have access to direct advice from the QIIME 2 development team at the
QIIME 2 forum. This helps researchers in my group make informed decisions, and the forum contributes to a
350 continuous development and improvement of the pipeline. Also available are tutorials that clarify how the
different plugins work. Without a background in bioinformatics and statistics, the forum and tutorials are an
352 absolute must to be able to run the plugins correctly. All in all, QIIME 2 is a free tool making complex and
advanced data curation possible for me as a researcher with no pre-existing knowledge of bioinformatics or
354 computer programming. From my experience, it has made it possible for researchers in my research group
to obtain hands-on knowledge with this new and complex field of research, rather than leaving bioinformatics
356 to externals.

Pedro J. Torres

Ph.D. Candidate

Department of Biology, San Diego State University, San Diego, California

360 Mr. Torres is a co-author of this paper.

362 The ability to add on plugins as needed and seeing the transition of QIIME from a platform for amplicon
based (specifically 16S rRNA) to multiomics analysis is very exciting and needed in a time when using
different omics technologies in a single study is starting to become the norm.

364 One of the greatest issues I have seen time and time again in analysis of next generation sequencing (NGS)
data is, as with many other things, REPRODUCIBILITY! Often when analyzing large amounts of data, good
366 notes for reproducibility later is not on the top of everyone's list. QIIME 2's automated data provenance
allows for tracking of not only the parameters and versions used to generate your data, but also the machine
368 used allows for transparent microbiome research without the need to even think about it. This is a big step
forward in allowing for data transparency and reproducibility.

370 I have also used many different bioinformatics software tools and it can be quite frustrating learning about a
new tool, wanting to use a new tool only to find out that there are issues that keep piling up and the authors
372 do not respond with help. The QIIME community is amazing! Any issue that I have had has been answered
in an appropriate time without the sarcasm or bullying you can encounter in other community forums. QIIME
374 2 has really done a great job at creating not only a tool but an online community resource well suited for the
novice and experienced data scientist alike unlike any I have ever seen.

376 Some of the interesting new plugins that are helping our lab advance our research is q2-longitudinal to look
378 at factors influencing alpha and beta diversity changes over time and q2-metabolomics to import our
metabolomics data into QIIME 2 for analysis.

Jonathan Warren

380 **National Laboratory Service, Environment Agency, Starcross, UK**

Mr. Warren is a co-author of this paper.

382 I started using QIIME 2 as an introduction to bioinformatics with no previous experience. It made learning the
384 basics of manipulating my DNA data much easier, and allowed me to try out the packaged tools much easier
than if they were to be used standalone. I like the fact QIIME 2 is open source as I believe that it is better for
386 the community and makes using DNA sequence data much easier for those without a computer science
background, as it can be tried for free with no cost to the user.

388 I contributed to QIIME 2 as I wanted to just import a lot of fastq files at once and I knew they were already in
a very specific format, and so I started writing some code in order to import using a whole folder at once.
Once I did this and tested it, I shared my code for all to use.

390 References

- 392 1. Koskinen, K. et al. *mBio* **8** (6), e00824-17 (2017). doi: 10.1128/mBio.00824-17
2. Schwendner, P. et al. *Microbiome* **5** (1), 129. doi:10.1186/s40168-017-0345-8
3. Mahnert, A. et al. *Front. Microbiol.* (2018). doi: 10.3389/fmicb.2018.02985
- 394 4. Caporaso, J.G. et al. *Nature Methods* (2010). Doi: 10.1038/nmeth.f.303
5. Bai, M. et al. *Microb. Ecol.* (2018). doi.:10.1007/s00248-018-1235-8.
- 396 6. Xie, N. et al. *Env. Microbiol.* 20(8), 3042-3056 (2018).
7. Tangedal et al., *Respir. Res.* 18, 164 (2017).