

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 Abnet²,
Gabriel A Al-Ghalith³, Harriet Alexander^{4,5}, Eric J Alm^{6,7}, Manimozhiyan Arumugam⁸, Francesco
6 Asnicar⁹, Yang Bai^{10,11,12}, Jordan E Bisanz¹³, Kyle Bittinger^{14,15}, Asker Brejnrod¹⁶, Colin J
Brislawn¹⁷, C Titus Brown⁵, Benjamin J Callahan^{18,19}, Andrés Mauricio Caraballo-Rodríguez²⁰,
8 John Chase¹, Emily Cope^{1,21}, Ricardo Da Silva²⁰, Pieter C Dorrestein²⁰, Gavin M Douglas²²,
Daniel M Durall²³, Claire Duvallet⁶, Christian F Edwardson²⁴, Madeleine Ernst²⁰, Mehrbod
10 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
12 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,
14 Dan Knights^{42,32}, Irina Koester^{43,20}, Tomasz Kosciolk⁴⁴, Jorden Kreps¹, Morgan GI Langille⁴⁵,
Joslynn Lee⁴⁶, Ruth Ley^{47,48}, Yong-Xin Liu^{10,11}, Erikka Lofffield², Catherine Lozupone⁴⁹, Massoud
16 Maher⁵⁰, Clarisse Marotz³⁰, Bryan Martin⁵¹, Daniel McDonald³⁰, Lauren J McIver^{35,36}, Alexey V
Melnik²⁰, Jessica L Metcalf⁵², Sydney C Morgan⁵³, Jamie T Morton^{30,50}, Ahmad Turan Naimey¹,
18 Jose A Navas-Molina^{50,30,54}, Louis Felix Nothias²⁰, Stephanie B Orchanian⁵⁵, Talima Pearson¹,
Samuel L Peoples^{56,57}, Daniel Petras²⁰, Mary Lai Preuss⁵⁸, Elmar Priesse⁴⁹, Lasse Buur
20 Rasmussen¹⁶, Adam Rivers⁵⁹, Michael S Robeson, II⁶⁰, Patrick Rosenthal⁵⁸, Nicola Segata⁹,
Michael Shaffer^{49,61}, Arron Shiffer¹, Rashmi Sinha², Se Jin Song³⁰, John R Spear⁶², Austin D
22 Swafford⁵⁵, Luke R Thompson^{63,64}, Pedro J Torres⁶⁵, Pauline Trinh⁶⁶, Anupriya Tripathi^{20,30,67},
Peter J Turnbaugh⁶⁸, Sabah Ul-Hasan⁶⁹, Justin JJ van der Hooft⁷⁰, Fernando Vargas⁶⁷, Yoshiki
24 Vázquez-Baeza³⁰, Emily Vogtmann², Max von Hippel⁷¹, William Walters⁴⁷, Yunhu Wan², Mingxun
Wang²⁰, Jonathan Warren⁷², Kyle C Weber^{59,73}, Chase HD Williamson¹, Amy D Willis⁷⁴,
26 Zhenjiang Zech Xu³⁰, Jesse R Zaneveld⁷⁵, Yilong Zhang⁷⁶, Rob Knight^{30,77,55}, and J Gregory
Caporaso^{1,21,+}

28 * These authors contributed equally to this work.

+ Please address correspondence about this document to gregcaporaso@gmail.com.

30 Author affiliations are provided following the Main Text References.

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

32

Abstract

34 We present QIIME 2, an open-source microbiome data science platform accessible to users
36 spanning the microbiome research ecosystem, from scientists and engineers to clinicians and
38 policy makers. QIIME 2 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 science.

Main text

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

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

58 QIIME 2 is developed based on a plugin architecture (Figure S1) that allows third-parties to
60 contribute functionality (see <https://library.qiime2.org>). QIIME 2 plugins exist for latest-generation
62 tools for sequence quality control from different sequencing platforms (DADA2¹⁰ and Deblur¹¹),
64 taxonomy assignment¹², and phylogenetic insertion¹³, that quantitatively improve results over
QIIME 1 and other tools (detailed in the corresponding tool-specific publications). Plugins also
support qualitatively new functionality including microbiome paired-sample and time-series
analysis¹⁴, critical for studying the impact of treatment on the microbiome, and for machine
learning¹⁵, including the ability to save trained models and apply them to new data and to
interrogate models to identify important microbiome features. Several recently released plugins,
66 including q2-cscs¹⁶, q2-metabolomics¹⁷, q2-shogun¹⁸, q2-metaphlan2¹⁹, and q2-picrust2²⁰,
provide initial support for analysis of metabolomics and shotgun metagenomics data. This marks

68 the potential of QIIME 2 to serve not only as a marker gene analysis tool, but also a
70 multi-dimensional and powerful data science platform that can be rapidly adapted to analyze
diverse microbiome features.

QIIME 2 provides many new interactive visualization tools (Figure 1), facilitating exploratory
72 analyses and result reporting. QIIME 2 View (<https://view.qiime2.org>) is a unique new service
(see Online Methods) that allows users to securely share and interact with results without
74 installing QIIME 2. The QIIME 2 visualizations presented in Figure 1 are provided in
Supplementary File 1 for readers to interact with using QIIME 2 View.

76 Reproducibility, transparency, and clarity of microbiome data science are guiding principles in
the QIIME 2 design. Toward this end, it includes a decentralized data provenance tracking
78 system: details of all analysis steps with references to intermediate data are automatically
stored in the results. Users can thus retrospectively determine exactly how any result was
80 generated (Figure 2). QIIME 2 also detects corrupted results, indicating that provenance is no
longer reliable and the results no longer contain information enabling reproducibility.
82 Provenance of the visualizations presented in Figure 1 can be interactively reviewed by loading
the contents of Supplementary File 1 with QIIME 2 View, providing far more detailed information
84 than can typically be provided in Methods text. QIIME 2 results are additionally semantically
typed (Figure 2) and actions indicate acceptable input types, clarifying the data that actions
86 should be applied to and making complex workflows less error-prone.

Finally, QIIME 2 provides a software development kit (see <https://dev.qiime2.org>) that can be
88 used to integrate it as a component of other systems (e.g., such as Qiita²¹ or Illumina
BaseSpace) and to develop interfaces targeted toward users with different levels of
90 computational sophistication (Figure S2). QIIME 2 provides the QIIME 2 Studio graphical user
interface and QIIME 2 View, interfaces designed for end-user biologists, clinicians, and policy
92 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,
94 providing a command line interface and Common Workflow Language²² wrappers for QIIME 2,
designed for high-performance computing experts.

96 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
98 microbiome data science.

Figures and figure captions

100 **Figure 1:** QIIME 2 provides many interactive visualization tools. Interactive versions of these
102 screen captures are available in Supplementary File 1 and at <https://github.com/qiime2/paper1>.
104 Detailed descriptions and methods are included in Online Methods. (A) Unweighted UniFrac
106 PCoA plot containing 37,680 samples, illustrating the scalability of QIIME 2. Colors indicate
108 sample type as described by the Earth Microbiome Project ontology (EMPO). (B) A feature
110 volatility plot illustrating change in *Bifidobacterium* abundance over time in breast-fed and
112 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.

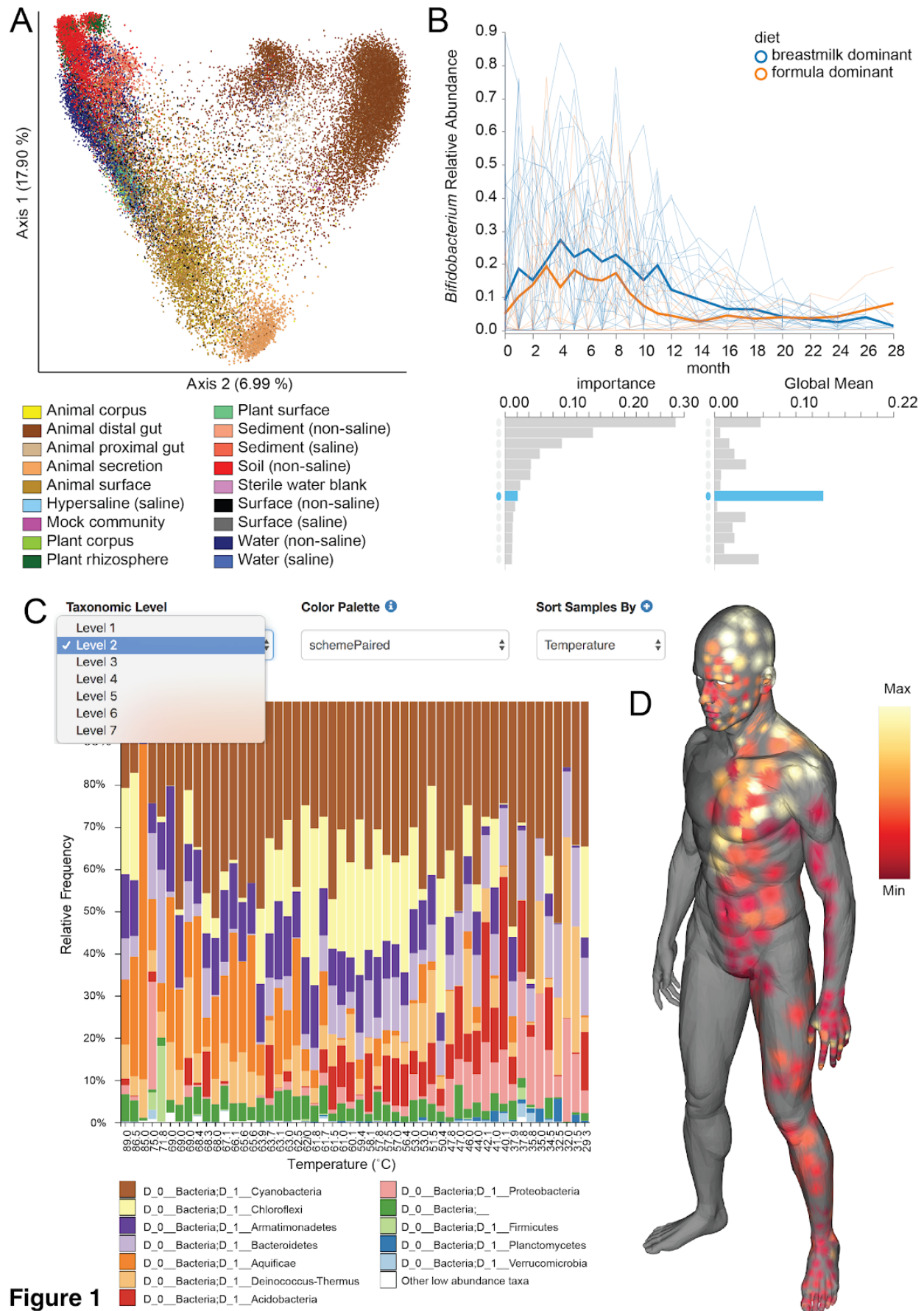


Figure 1

114 **Figure 2:** QIIME 2 iteratively records data provenance, ensuring bioinformatics reproducibility.
116 This simplified diagram illustrates the automatically tracked information about the creation of the
118 taxonomy barplot presented in Figure 1c. QIIME 2 results (circles) contain network diagrams
120 illustrating the data provenance stored in the result. Actions (quadrilaterals) are applied to
122 QIIME 2 results and generate new results. Arrows indicate flow of QIIME 2 results through
124 actions. TaxonomicClassifier and FeatureData[Sequence] inputs contain independent
126 provenance (red and blue, respectively) and are provided to a classify action (yellow), which
128 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.

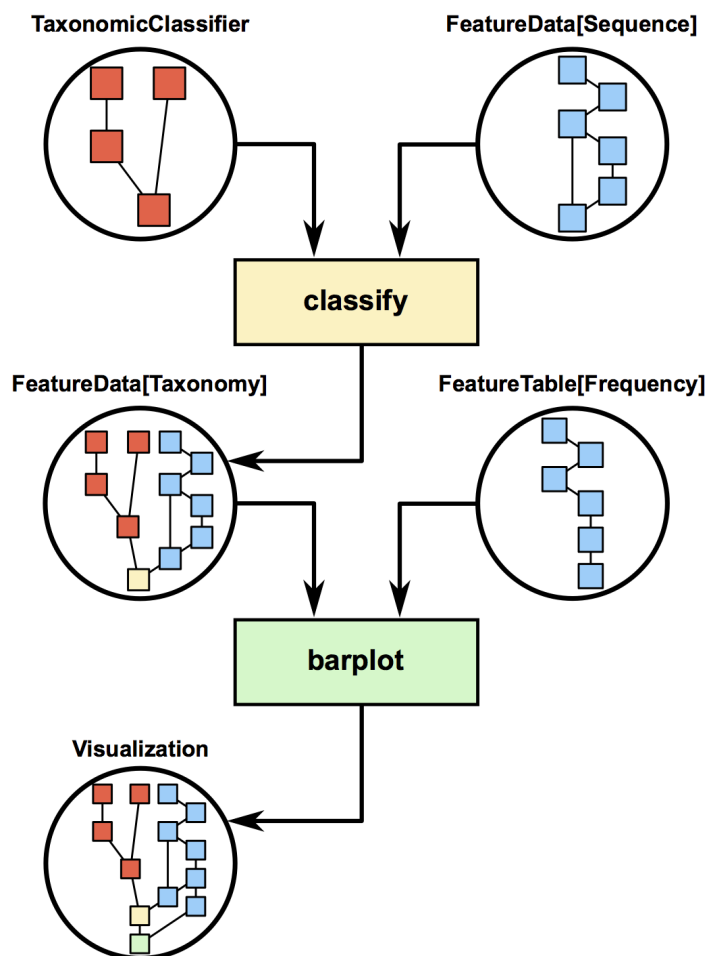


Figure 2

Code availability

130 QIIME 2 is open source and free for all use, including commercial. It is licensed under the BSD
132 3-clause license. Source code is available at <https://github.com/qiime2>.

Data availability

Data for the analyses presented in Figure 1 are available as follows: (a) Earth Microbiome
134 Project data was obtained from <ftp://ftp.microbio.me/emp/release1>, and the American Gut
136 Project (AGP) data was obtained from Qiita (<http://qiita.microbio.me>) study 10317. (b) Sequence
138 data are available in Qiita under study id 10249 and EBI under accession number ERP016173.
140 (c) Sequence data are available in Qiita under study id 925 and EBI 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 <https://github.com/qiime2/paper1>.

Acknowledgements

142 QIIME 2 development was primarily funded by NSF Award 1565100 to JGC and RK. Partial
support was also provided from the following grants: NIH U54CA143925 (JGC, TP) and
144 U54MD012388 (JGC, 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
146 QYZDB-SSW-SMC021 (YB); from the Australian National Health and Medical Research Council
APP1085372 (GAH, JGC, Von Bing Yap and RK); and from Natural Sciences and Engineering
148 Research Council (NSERC) to DLG. Thanks to the Yellowstone Center for Resources for
research permit #5664 to JRS for Yellowstone access and sample collection. We would like to
150 thank the users of QIIME 1 and 2, whose invaluable feedback has shaped QIIME 2.

Author contributions

152 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
154 JGC developed 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
156 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
158 manuscript and/or testing of the QIIME 2. CA, 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,
160 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,

162 PT, AT, JJJV, YV, MV, MW, KCW, ADW, ZZX, JRZ, YZ, and JGC contributed software to QIIME
164 2 plugins, interfaces, framework, and/or build and test systems.

Main text references

1. Smith, M.I. et al. *Science* **339**, 548–554 (2013).
- 166 2. Gopalakrishnan, V. et al. *Science* **359**, 97–103 (2018).
- 168 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).
- 170 4. Lee, K., Pletcher, S.D., Lynch, S.V., Goldberg, A.N. & Cope, E.K. *Front. Cell. Infect. Microbiol.* **8**, 168 (2018).
- 172 5. Metcalf, J.L. et al. *Science* **351**, 158–162 (2016).
- 174 6. Rubin, R.L. et al. *Ecol. Appl.* **28**, 1594–1605 (2018).
- 176 7. Pineda, A., Kaplan, I. & Bezemer, T.M. *Trends Plant Sci.* **22**, 770–778 (2017).
- 178 8. Kaponó, C.A. et al. *Sci. Rep.* **8**, 3669 (2018).
- 180 9. Barr, T. et al. *Gut Microbes* 1–44 (2018).
- 182 10. Callahan, B.J. et al. *Nat. Methods* (2016).doi:10.1038/nmeth.3869
- 184 11. Amir, A. et al. *mSystems* **2**, (2017).
- 186 12. Bokulich, N.A. et al. *Microbiome* **6**, 90 (2018).
- 188 13. Janssen, S. et al. *mSystems* **3**, e00021–18 (2018).
14. Bokulich, N. et al. *bioRxiv* 223974 (2017).doi:10.1101/223974
15. Bokulich, N. et al. *bioRxiv* 306167 (2018).doi:10.1101/306167
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

Author affiliations

190 ¹Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA. ²Metabolic
192 Epidemiology Branch, National Cancer Institute, Rockville, MD, USA. ³Department of Computer
194 Science and Engineering, University of Minnesota, Minneapolis, Minnesota, USA. ⁴Biology
196 Department, Woods Hole Oceanographic Institution, Woods Hole, MA, USA. ⁵Department of
198 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
200 Genomics, Institute of Genetics and Developmental Biology, Chinese Academy of Sciences,
Beijing, China. ¹¹Centre of Excellence for Plant and Microbial Sciences (CEPAMS), Institute of

202 Genetics and Developmental Biology, Chinese Academy of Sciences & John Innes Centre,
Beijing, China. ¹²University of Chinese Academy of Sciences, Beijing, China. ¹³Department of
204 Microbiology and Immunology, University of California, San Francisco, CA, USA. ¹⁴Division of
Gastroenterology and Nutrition, Children's Hospital of Philadelphia, Philadelphia, PA, USA.
206 ¹⁵Hepatology, Children's Hospital of Philadelphia, Philadelphia, PA, USA. ¹⁶Novo Nordisk
Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences,
208 University of Copenhagen, Denmark. ¹⁷Earth and Biological Sciences Directorate, Pacific
Northwest National Laboratory, Richland, WA, USA. ¹⁸Department of Population Health &
210 Pathobiology, North Carolina State University, Raleigh, NC, USA. ¹⁹Bioinformatics Research
Center, North Carolina State University, Raleigh, NC, USA. ²⁰Collaborative mass spectrometry
212 innovation center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of
California San Diego, San Diego, CA, USA. ²¹Department of Biological Sciences, Northern
214 Arizona University, Flagstaff, AZ, USA. ²²Department of Microbiology and Immunology,
Dalhousie University, Halifax, Nova Scotia, Canada. ²³Irving K. Barber School of Arts and
216 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
218 Aquarium, Chicago, IL, USA. ²⁵Department of Biology, University of British Columbia Okanagan,
Okanagan, BC, Canada. ²⁶Computational Bioscience Graduate Program, University of Colorado
220 Denver Anschutz Medical Campus, Aurora, Colorado, USA. ²⁷Department of Medicine, Division
of Biomedical Informatics and Personalized Medicine, University of Colorado Denver Anschutz
222 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
224 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,
226 Michigan State University, East Lansing, MI, USA. ³²Department of Computer Science and
Engineering, University of Minnesota, Minneapolis, MN, USA. ³³Stanford University, Statistics
228 Department, Palo Alto, CA, USA. ³⁴Department of Computer Science and Engineering,
University of California San Diego, La Jolla, CA, USA. ³⁵Department of Biostatistics, Harvard
230 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,
232 Canberra, ACT, Australia. ³⁸Department of Pediatric Oncology, Hematology and Clinical
Immunology, Heinrich-Heine University Dusseldorf, Dusseldorf, Germany. ³⁹Department of
234 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
236 State University, Department of Biology, San Diego, CA, USA. ⁴²Biotechnology Institute,
University of Minnesota, Saint Paul, MN, USA. ⁴³Scripps Institution of Oceanography, University
238 of California San Diego, La Jolla, CA, USA. ⁴⁴Department of Pediatrics, University of California
San Diego, La Jolla, CA. ⁴⁵Department of Pharmacology, Dalhousie University, Halifax, Nova
240 Scotia, Canada. ⁴⁶Science Education, Howard Hughes Medical Institute, Ashburn, VA, USA.
⁴⁷Department of Microbiome Science, Max Planck Institute for Developmental Biology,
242 Tübingen, Germany. ⁴⁸Department of Molecular Biology and Genetics, Cornell University, Ithaca,
NY, USA. ⁴⁹Department of Medicine, Division of Biomedical Informatics and Personalized
244 Medicine, University of Colorado Denver Anschutz Medical Campus, Aurora, CO, USA.

246 ⁵⁰Department of Computer Science & Engineering, University of California San Diego, La Jolla,
CA, USA. ⁵¹Department of Statistics, University of Washington, Seattle, WA, USA. ⁵²Department
248 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.
250 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
252 University, Syracuse, NY, USA. ⁵⁷School of STEM, University of Washington Bothell, Bothell,
WA, USA. ⁵⁸Department of Biological Sciences, Webster University, St Louis, MO, USA.
254 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.
256 ⁶¹Computational Bioscience Program, University of Colorado Denver Anschutz Medical
Campus, Aurora, CO, USA. ⁶²Department of Civil and Environmental Engineering, Colorado
258 School of Mines, Golden, CO, USA. ⁶³Department of Biological Sciences and Northern Gulf
Institute, University of Southern Mississippi, Hattiesburg, Mississippi, USA. ⁶⁴Ocean Chemistry
260 and Ecosystems Division, Atlantic Oceanographic and Meteorological Laboratory, National
Oceanic and Atmospheric Administration, La Jolla, CA, USA. ⁶⁵Department of Biology, San
262 Diego State University, San Diego, CA, USA. ⁶⁶Department of Environmental and Occupational
Health Sciences, University of Washington, Seattle, WA, USA. ⁶⁷Division of Biological Sciences,
264 University of California San Diego, San Diego, CA, USA. ⁶⁸Department of Microbiology and
Immunology, University of California San Francisco, San Francisco, CA, USA. ⁶⁹Quantitative
266 and Systems Biology Graduate Program, University of California Merced, Merced, CA, USA.
⁷⁰Bioinformatics Group, Wageningen University, Wageningen, The Netherlands. ⁷¹Department of
268 Mathematics, University of Arizona, Tucson, AZ, USA. ⁷²National Laboratory Service,
Environment Agency, Starcross, UK. ⁷³College of Agriculture and Life Sciences, University of
270 Florida, Gainesville, FL, USA. ⁷⁴Department of Biostatistics, University of Washington, Seattle,
WA, USA. ⁷⁵University of Washington Bothell, School of STEM, Division of Biological Sciences,
272 Bothell, WA, USA. ⁷⁶Merck & Co. Inc., Kenilworth, NJ, USA. ⁷⁷Department of Computer Science
and Engineering, University of California San Diego, La Jolla, California, USA

274 QIIME 2: Reproducible, interactive, 276 scalable, and extensible microbiome data science

Online methods

278 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>.

282 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
284 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
286 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
288 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
290 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
292 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>).
294 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
296 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
298 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
300 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.

302 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
304 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
306 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
308 unique identifiers (UUIDs) to each execution of an **Action**, and to all **Results**.

310 QIIME 2 stores information about the series of Actions that led to a Result, along with
312 information about the environment (including versions of all QIIME 2 packages and other Python
314 dependencies) where each Action was executed, and the data itself. We refer to this process as
316 data provenance tracking, or simply provenance tracking. We did not want to create new
318 bioinformatics file formats to support the storage of data provenance, so QIIME 2 Results are
320 instead stored as zip files containing a data directory that contains only the data in a relevant
322 format (e.g., fasta or fastq for sequence data, newick for phylogenetic trees, etc), plus
324 QIIME-2-specific metadata in other directories (such as provenance). These files use the
326 extension .qza (for QIIME zipped artifact) or .qzv (for QIIME zipped visualization), but they are
328 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.

330 Another important component of QIIME 2 is its **semantic type system**. All Artifacts
332 used in QIIME 2 are annotated with a semantic description of their type which conveys the
334 meaning of the data. Semantic types differ from data types (how data is represented in memory)
336 or file formats (how data is stored on disk), and allow QIIME 2 to constrain the composition of
338 multiple actions to only those combinations which are semantically meaningful without needing
340 to consider the specific file formats or data types. This also makes it possible to determine what
342 Actions could be applied (and in what order) to generate a given Artifact from some set of input
344 Artifacts. For example, phylogenetic trees in QIIME 2 can be either rooted or unrooted, and
346 these two concepts are represented by the semantic types `Phylogeny[Rooted]` and
`Phylogeny[Unrooted]`, respectively. QIIME 2 could support loading these into multiple
348 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.

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

354

QIIME 2 View

356 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
358 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
360 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
362 `.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
364 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
366 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
368 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
370 allows QIIME 2 View to generate a dynamic provenance visualization, summarizing the entire
provenance of the archive in question.

372

Extracting data from QIIME 2 archives

374 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`
376 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
378 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,
380 users could follow these steps to access QIIME 2 Results.

382 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
```

384 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
386 `8dc793b8-7284-462a-8578-6370ffccebd`.

```

388 $ unzip rep-seqs.qza
Archive: rep-seqs.qza
   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/metadata.yaml
390   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/VERSION
   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/metadata.yaml
392   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/citations.bib
   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/VERSION
394   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/me
396 tadata.yaml
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/ci
398 tations.bib
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/VE
400 RSION
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/bdaa3214-f883-4c8b-8db3-f6ea4910d724/ac
402 tion/action.yaml
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/me
406 tadata.yaml
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/ci
408 tations.bib
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/VE
412 RSION
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/ac
414 tion/action.yaml
   inflating:
8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/artifacts/7097fc98-ad5f-4b9d-a33e-39cd36857a0d/ac
416 tion/barcodes.tsv
   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/provenance/action/action.yaml
420   inflating: 8dc793b8-7284-462a-8578-6370ffcceebdc/data/dna-sequences.fasta
422

```

424 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
426 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
428 file containing sequence identifiers and sequences. The first four lines of this file can be viewed
as follows:

```

430 $ head -4 8dc793b8-7284-462a-8578-6370ffcceebdc/data/dna-sequences.fasta
>f352c1f1efecf483511c2270aab0ae6
TACGTAGGGTGCAGCGTTAATCGGAATTACTGGCGTAAAGCGTGCCGAGCGGTTTTGTGTAAGACAGAGGTGAAATCCCCGGGCT
432 CAACCTGGGAAGTGCCTTTGTGACTGCAAGGCTG
>82e72255267397b777a1afd44ea22755
TACGGAGGATCCAAGCGTTATCCGGAATCATTGGGTTTAAAGGGTCCGTAGCGGTTTAGTAAGTCAGTGGTAAAAGCCCATCGCT
434 CAACGGTGGAAACGGCCATTGATACTGCTAGACTT

```

436

QIIME 2 user and developer community

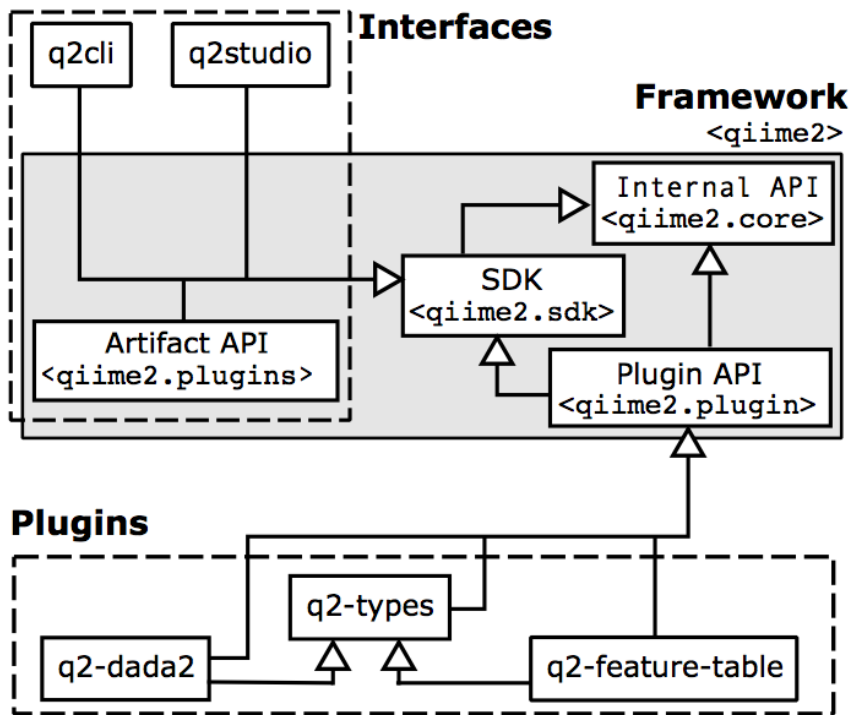
438 QIIME 2 officially succeeded QIIME 1 (<http://www.qiime.org>)³ in January of 2018, and has
440 developed an engaged user base and community. As of this writing there are over 1980 active
442 users (users who have performed an action, such as creating or liking a post) on the QIIME 2
444 Forum; over 3000 monthly downloads of QIIME 2 from Anaconda; over 8000 unique visitors to
446 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.

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

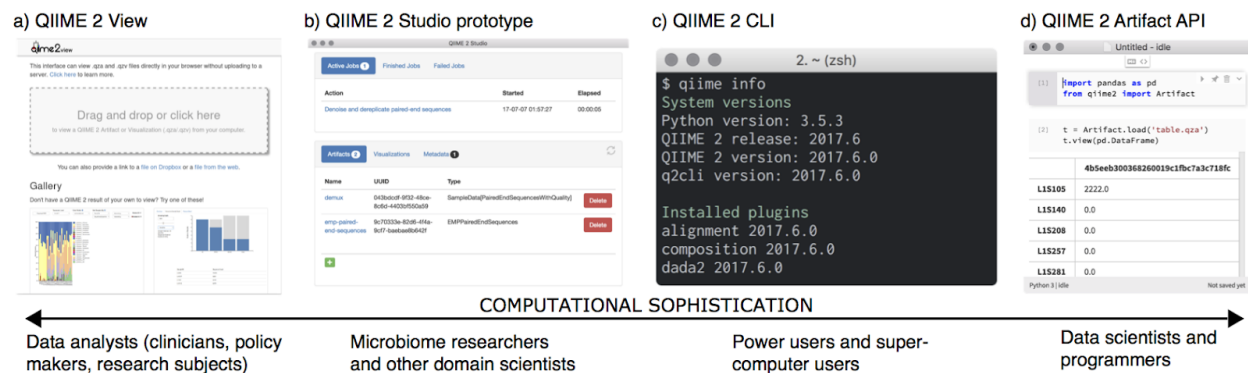
450

Supplementary figures, files, and captions

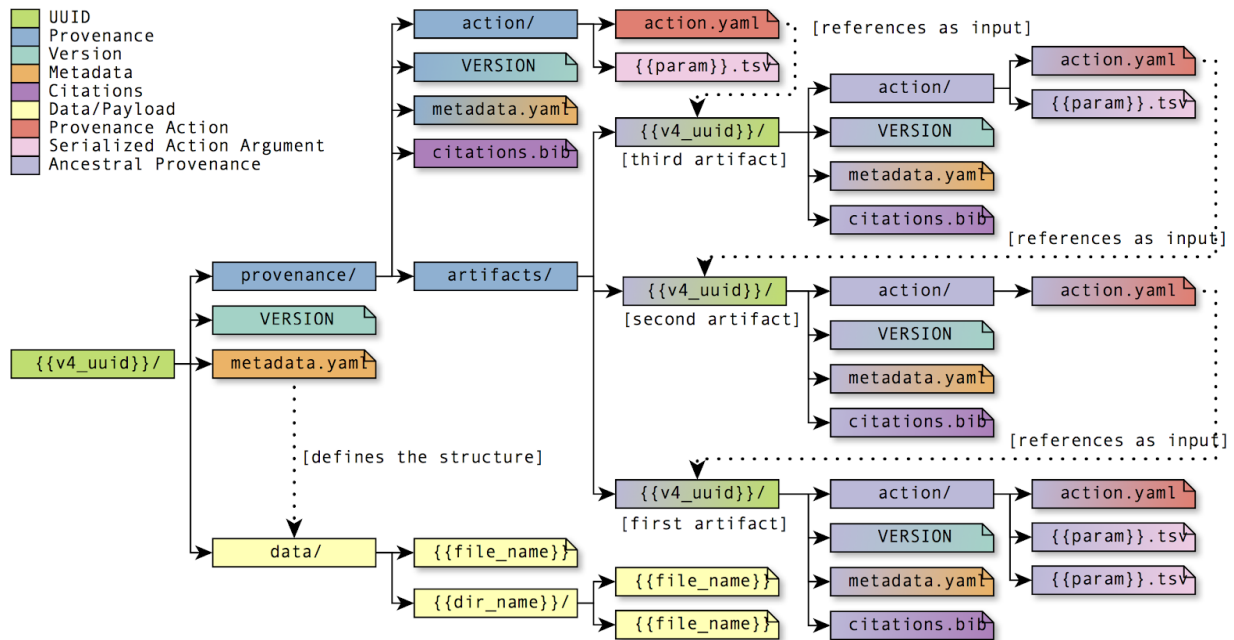
452 **Figure S1. Schematic diagram of the QIIME 2 system.** Interfaces define how users interact
with the system; plugins define all domain-specific functionality; and the framework mediates
454 communication between plugins and interfaces, and performs core functionality such as
provenance tracking. Arrows indicate dependencies. Interfaces interact only with the
456 `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.



458 **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
 460 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
 462 convenient for lead investigators, clinicians, or policy makers who may want to explore
 interactive visualizations generated by others. (b) Researchers who prefer graphical interfaces
 464 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
 466 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
 468 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”.



470 **Figure S3. Anatomy of a QIIME 2 Archive (i.e., .qza or .qzv file).** QIIME 2 stores data in a
 472 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.



474 **Supplementary File 1** contains the QIIME 2 `.qzv` files corresponding to **Figure 1a-d**. These
475 are also accessible at <https://github.com/qiime2/paper1> and can be viewed using QIIME 2 View
476 (<https://view.qiime2.org>) where readers can interact with the results, and explore the methods
477 used to generate them (see the Provenance tab after loading a `.qzv` file with QIIME 2 View).
478 We describe the methods used to generate each of these visualizations here, and this
479 information can be compared to the data provenance which contains far more detail that is
480 possible or desirable to include in supplementary methods text.

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

b-feature-volatility.qzv: Data were generated on five sequencing runs of V4 16S rRNA gene
490 amplicons from the ECAM study¹⁰. Forward reads were imported separately in
491 `EMPSingleEndDirFmt` format, demultiplexed with `q2-demux's emp_single` method, and
492 denoised using `q2-dada2's denoise_single` method (`trunc_len=150`, other parameters
493 used default values)¹¹. Denoised feature tables and sequences were merged using
494 `q2-feature-table's merge` and `merge-seqs` methods, respectively. `q2-feature-table's`
495 `filter-samples` method was used to remove samples with fewer than 2000 sequences, and
496 to perform metadata-based filtering to retain only children's samples. A naive Bayes taxonomy
497 classifier was trained on the Greengenes¹² reference sequences (clustered at 99% similarity)
498 using `q2-feature-classifier's fit-classifier-naive-bayes` method¹³. This classifier was
499 used to taxonomically classify the ECAM ASVs using `q2-feature-classifier's`
500 `classify-sklearn` method¹³. ASVs were collapsed based on genus-level taxonomy using
501 `q2-taxa's collapse` method. Temporally predictive features were identified using `q2-longitudinal's`
502 `feature-volatility` pipeline¹⁴ using default parameters. Data contained in this artifact have
503 been described in a previous publication¹⁴.

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

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

516

Online methods references

1. Grüning, B. et al. *Nat. Methods* **15**, 475–476 (2018).
- 518 2. Huber, W. et al. *Nat. Methods* **12**, 115–121 (2015).
3. Caporaso, J.G. et al. *Nat. Methods* **7**, 335–336 (2010).
- 520 4. Thompson, L.R. et al. *Nature* **551**, 457–463 (2017).
5. McDonald, D. et al. *mSystems* **3**, e00031–18 (2018).
- 522 6. McDonald, D. et al. *Gigascience* **1**, 7 (2012).
7. Amir, A. et al. *mSystems* **2**, (2017).
- 524 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).
- 526 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
- 528 12. McDonald, D. et al. *ISME J.* **6**, 610–618 (2012).
13. Bokulich, N.A. et al. *Microbiome* **6**, 90 (2018).
- 530 14. Bokulich, N. et al. *bioRxiv* 223974 (2017).doi:10.1101/223974
15. Protsyuk, I. et al. *Nat. Protoc.* **13**, 134–154 (2018).
- 532 16. Bouslimani, A. et al. *Proc. Natl. Acad. Sci. U. S. A.* **112**, E2120–9 (2015).