

 Open access • Posted Content • DOI:10.1101/2021.06.04.446988

## Multi-omics profiling of Earth's biomes reveals that microbial and metabolite composition are shaped by the environment — [Source link](#)

Justin P. Shaffer, Louis-Félix Nothias, Louis-Félix Nothias, Luke R. Thompson ...+55 more authors

**Institutions:** University of California, San Diego, University of Montana, Mississippi State University, Atlantic Oceanographic and Meteorological Laboratory ...+11 more institutions

**Published on:** 06 Jun 2021 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Metabolome and Metagenomics

Related papers:

- [Functional analysis of pristine estuarine marine sediments](#)
- [Filamentous fungi from extreme environments as a promising source of novel bioactive secondary metabolites.](#)
- [Predicted Relative Metabolomic Turnover \(PRMT\): determining metabolic turnover from a coastal marine metagenomic dataset](#)
- [Web of microbes \(WoM\): a curated microbial exometabolomics database for linking chemistry and microbes.](#)
- [Linking soil biology and chemistry in biological soil crust using isolate exometabolomics.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/multi-omics-profiling-of-earth-s-biomes-reveals-that-1c25cz093r>

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

Multi-omics profiling of Earth's biomes reveals that microbial and metabolite composition are shaped by the environment

Justin P. Shaffer<sup>1,#</sup>, Louis-Félix Nothias<sup>2,3,#</sup>, Luke R. Thompson<sup>4,5,#</sup>, Jon G. Sanders<sup>6</sup>, Rodolfo A. Salido<sup>7</sup>, Sneha P. Couvillion<sup>8</sup>, Asker D. Brejnrod<sup>3</sup>, Shi Huang<sup>1,9</sup>, Franck Lejzerowicz<sup>1,9</sup>, Holly L. Lutz<sup>1,10</sup>, Qiyun Zhu<sup>11,12</sup>, Cameron Martino<sup>9,13</sup>, James T. Morton<sup>14</sup>, Smruthi Karthikeyan<sup>1</sup>, Mélissa Nothias-Esposito<sup>2,3</sup>, Kai Dührkop<sup>15</sup>, Sebastian Böcker<sup>15</sup>, Hyunwoo Kim<sup>10</sup>, Alexander A. Aksenov<sup>2,3</sup>, Wout Bittremieux<sup>2,3,16</sup>, Jeremiah J. Minich<sup>10</sup>, Clarisse Marotz<sup>1</sup>, MacKenzie M. Bryant<sup>1</sup>, Karenina Sanders<sup>1</sup>, Tara Schwartz<sup>1</sup>, Greg Humphrey<sup>1</sup>, Yoshiki Vásquez-Baeza<sup>9</sup>, Anupriya Tripathi<sup>1,3</sup>, Laxmi Parida<sup>17</sup>, Anna Paola Carrieri<sup>18</sup>, Niina Haiminen<sup>17</sup>, Kristen L. Beck<sup>19</sup>, Promi Das<sup>1,10</sup>, Antonio González<sup>1</sup>, Daniel McDonald<sup>1</sup>, Søren M. Karst<sup>20</sup>, Mads Albertsen<sup>21</sup>, Gail Ackermann<sup>1</sup>, Jeff DeReus<sup>1</sup>, Torsten Thomas<sup>22</sup>, Daniel Petras<sup>2,10,23</sup>, Ashley Shade<sup>24</sup>, James Stegen<sup>8</sup>, Se Jin Song<sup>9</sup>, Thomas O. Metz<sup>8</sup>, Austin D. Swafford<sup>9</sup>, Pieter C. Dorrestein<sup>2,3</sup>, Janet K. Jansson<sup>8</sup>, Jack A. Gilbert<sup>1,10</sup>, Rob Knight<sup>1,7,9,25,\*</sup>, and the Earth Microbiome Project 500 (EMP500) Consortium

<sup>1</sup>Department of Pediatrics, School of Medicine, University of California San Diego, La Jolla, California, USA.

<sup>2</sup>Collaborative Mass Spectrometry Innovation Center; University of California San Diego; La Jolla, CA 92093; USA

<sup>3</sup>Skaggs School of Pharmacy and Pharmaceutical Sciences; University of California San Diego; La Jolla CA 92093; USA

<sup>4</sup>Northern Gulf Institute, Mississippi State University, Mississippi State, Mississippi, USA

<sup>5</sup>Ocean Chemistry and Ecosystems Division, Atlantic Oceanographic and Meteorological Laboratory, National Oceanic and Atmospheric Administration, Miami, Florida, USA

<sup>6</sup>Department of Ecology and Evolutionary Biology, Cornell University, Ithaca, New York, USA.

<sup>7</sup>Department of Bioengineering, University of California San Diego, La Jolla, California, USA.

<sup>8</sup>Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, WA, USA

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

<sup>9</sup>Center for Microbiome Innovation, Jacobs School of Engineering, University of California San Diego, La Jolla, California, USA.

<sup>10</sup>Scripps Institution of Oceanography, University of California San Diego, La Jolla, California, USA.

<sup>11</sup>School of Life Sciences, Arizona State University, Tempe, AZ 85281, USA

<sup>12</sup>Biodesign Center for Fundamental and Applied Microbiomics, Arizona State University, Tempe, AZ 85281, USA

<sup>13</sup>Bioinformatics and Systems Biology Program, Jacobs School of Engineering, University of California San Diego, La Jolla, California, USA.

<sup>14</sup>Center for Computational Biology, Flatiron Institute, Simons Foundation

<sup>15</sup>Chair for Bioinformatics, Friedrich Schiller University, Jena, Germany

<sup>16</sup>Department of Computer Science, University of Antwerp, Antwerp, Belgium

<sup>17</sup>IBM Research, T.J. Watson Research Center, Yorktown Heights, NY, USA

<sup>18</sup>IBM Research Europe - Daresbury, UK

<sup>19</sup>IBM Research, Almaden Research Center, San Jose, CA, USA

<sup>20</sup>Department of Virus and Microbiological Special Diagnostics, Statens Serum Institute, Copenhagen, Denmark

<sup>21</sup>Department of Chemistry and Bioscience, Aalborg University, Aalborg, Denmark

<sup>22</sup>Centre for Marine Science and Innovation, School of Biological, Earth and Environmental Science, The University of New South Wales, Sydney, 2052, Australia

<sup>23</sup>Interfaculty Institute of Microbiology and Infection Medicine, University of Tübingen, Tübingen, Baden-Württemberg 72076, Germany

<sup>24</sup>Department of Microbiology and Molecular Genetics, Michigan State University East Lansing MI USA

<sup>25</sup>Department of Computer Science and Engineering, Jacobs School of Engineering, University of California San Diego, La Jolla, California, USA.

<sup>#</sup>Co-first author

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

**Abstract:** Microbes produce an array of secondary metabolites that perform diverse functions from communication to defense<sup>1</sup>. These metabolites have been used to benefit human health and sustainability<sup>2</sup>. In their analysis of the Genomes from Earth’s Microbiomes (GEM) catalog<sup>3</sup>, Nayfach and co-authors observed that, whereas genes coding for certain classes of secondary metabolites are limited or enriched in certain microbial taxa, “*specific chemistry is not limited or amplified by the environment, and that most classes of secondary metabolites can be found nearly anywhere*”. Although metagenome mining is a powerful way to annotate biosynthetic gene clusters (BCGs), chemical evidence is required to confirm the presence of metabolites and comprehensively address this fundamental hypothesis, as metagenomic data only identify metabolic potential. To describe the Earth’s metabolome, we use an integrated omics approach: the direct survey of metabolites associated with microbial communities spanning diverse environments using untargeted metabolomics coupled with metagenome analysis. We show, in contrast to Nayfach and co-authors, that the presence of certain classes of secondary metabolites can be limited or amplified by the environment. Importantly, our data indicate that considering the relative abundances of secondary metabolites (i.e., rather than only presence/absence) strengthens differences in metabolite profiles across environments, and that their richness and composition in any given sample do not directly reflect those of co-occurring microbial communities, but rather vary with the environment.

## Shaffer et al. Metabolite-microbe profiles are shaped by the environment

From a genomics perspective, information regarding metabolic potential is obtained through detection and classification of biosynthetic gene clusters (BGCs), the genomic loci underlying the production of secondary metabolites and their precursors<sup>4</sup>. The most sensitive approaches amplify BGC-specific marker sequences by PCR<sup>5</sup>, but only metagenomic methods can link BGCs to their genomes (i.e., metagenome-assembled genomes, or MAGs) of origin and detect BGCs in novel MAGs. Nayfach and co-authors uncovered 104,211 putative BGC regions from 52,515 microbial MAGs. Surprisingly, their analysis showed that, although the main classes of secondary metabolites are enriched in particular microbial taxa, the relative distribution of secondary metabolite biosynthetic potential across environments was conserved, implying that most classes of secondary metabolites are not “*limited or amplified*” by the environment. The authors acknowledged that most of their annotated BGCs had incomplete sequences, potentially impacting annotation and quantification, but that this was consistent with previous studies. More importantly, gene-level data about BGCs inferred from MAGs cannot offer information about actual synthesis (e.g., gene expression), creating uncertainty about the distribution of secondary metabolites across environments<sup>6-9</sup>. Even with high-coverage gene expression data, currently lacking for most environments, the complex structural and modular nature of many secondary metabolites prevents their accurate association with the underlying genomic origins<sup>10</sup>. Furthermore, quantifying metabolite diversity from such metatranscriptomic and/or metaproteomic data (also lacking for most environments) is problematic due to a suite of post-translational processes that can dissociate the level of gene transcription from the abundance of gene products<sup>11</sup>. Finally, shotgun metagenomics does not capture BGCs from low-abundance MAGs efficiently, as shown from comparative studies of targeted sequencing approaches<sup>5</sup>.

An approach to surmount these issues is to complement metagenomics with a direct survey of secondary metabolites using untargeted metabolomics. Liquid chromatography with untargeted tandem mass spectrometry (LC-MS/MS) is a versatile method that detects tens-of-thousands of metabolites in biological samples<sup>12</sup>. Although LC-MS/MS metabolomics has generally suffered from a low metabolite annotation rate when applied to non-model organisms, recent computational advances can systematically

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

classify metabolites using their fragmentation spectra<sup>13</sup>. Untargeted metabolomics provides the relative abundance (i.e., intensity) of each metabolite detected across samples rather than just counts of unique structures (e.g., Fig. 1a vs. 1b), and thus provides a direct readout of the surveyed environment, a result that is difficult to achieve with a purely genomics approach. While there is a clear need for the use of untargeted metabolomics to quantify the metabolic activities of microbiota, this methodology has been limited by the challenge of discriminating the secondary metabolites produced by microbes from tens-of-thousands of metabolites detected in the environment. To resolve this bottleneck, we devised a computational method for recognizing and annotating putative secondary metabolites of microbial origin from fragmentation spectra. The annotations were first obtained from spectral library matching and *in silico* annotation<sup>14</sup> using the GNPS web-platform<sup>15</sup>. These annotations were then queried against microbial metabolite reference databases (i.e., Natural Products Atlas<sup>16</sup> and MIBiG<sup>17</sup>), and molecular networking<sup>18</sup> was used to propagate the annotation to similar metabolites. Finally, a global chemical classification of these metabolites was achieved using SIRIUS/CANOPUS<sup>13</sup>.

We used this methodology to quantify microbial secondary metabolites from diverse microbial communities that span 20 major environments from the Earth Microbiome Project 500 (EMP500) dataset (Extended Data Fig. 1, Table S1). With this dataset, we show that although the presence/absence (i.e., occurrence) of major classes of microbially-related metabolites is indeed relatively conserved across habitats, their relative abundance reveals specific chemistry that is limited or amplified by the environment, especially at more resolved chemical class ontology levels (Fig. 1). Importantly, when considering differences in the relative abundances of all microbially-related metabolites, profiles among environments were so distinct that we could identify particular metabolites whose abundances were enriched in certain environments (Fig. 2a,c, Table S2, Table S3). For example, microbially-related metabolites associated with the carbohydrate pathway were especially enriched in aquatic samples, whereas those associated with the polyketide- and shikimate and phenylpropanoid pathways enriched in sediment, soil, and fungal samples (Fig. 2a). Interestingly, distinct analytical approaches confirmed specific metabolites as particularly important for distinguishing aquatic samples ( $C_{28}H_{58}O_{15}$ , pathway:

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

carbohydrates, superclass: glycerolipids), non-saline plant surface samples ( $C_{13}H_{10}O$ , pathway: shikimates and phenylpropanoids, superclass: flavonoids, class: chalcones), and non-saline animal distal gut samples ( $C_{24}H_{38}O_4$ , pathway: terpenoids, superclass: steroids, class: cholane steroids) (Figs. 2c, Tables S2, S3).

We also identified specific metabolites that could classify specific environments with 68.9% accuracy in machine-learning analysis (Fig. 3a, Extended Data Fig. 3, Table S4), and found further support for the importance of particular metabolites in distinguishing environments, including the putative cholane steroid above (i.e.,  $C_{24}H_{38}O_4$ ), and three metabolites enriched in non-saline soil and plant corpus samples (Fig. 2c, Fig. 3a, Table S4).

In addition to showing that the relative abundances of microbially-related metabolites distinguish environments, our results highlight the advantages of using a multi-omics approach to interpret and predict the contributions of microbes and their environments to chemical profiles in nature. Moreover, our approach illustrates that recent advances in computational annotation tools offer a powerful toolbox to interpret untargeted metabolomics data<sup>13</sup>. With these observations, we hypothesized that the differences in the relative abundances of particular metabolites among environments were due in part to underlying differences in microbial community composition and diversity. To begin to explore these relationships, we analyzed our shotgun metagenomics data and found strong correlations between microbially-related metabolite richness and microbial taxon richness for certain environments (i.e., Animal proximal gut (saline)  $r = 0.73$ ,  $p$ -value  $< 0.01$ ; Plant corpus (non-saline)  $r = 0.74$ ,  $p$ -value  $< 0.001$ ; Sediment (non-saline)  $r = 0.42$ ,  $p$ -value  $= 0.05$ ; Water (saline)  $r = 0.57$ ,  $p$ -value  $= 0.01$ ) (Fig. 2b; Table S5). We also found similarity in the clustering of samples by environment between datasets (Fig. 2c,d), and a strong correlation between sample-sample distances based on microbially related metabolites vs. microbial taxa (Extended Data Table 1). Using machine-learning, we determined that specific microbial taxa and their functions could classify environments with 72.08% and 68.19% accuracy, respectively (Extended Data Figs. 2 and 3). In addition, we examined correlations between microbe-metabolite co-occurrences learned from shotgun metagenomic profiles and (1) log-fold changes of metabolites across environments, and (2) global distributions of metabolites, and found strong relationships with each (Figure 3b). In particular, the

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

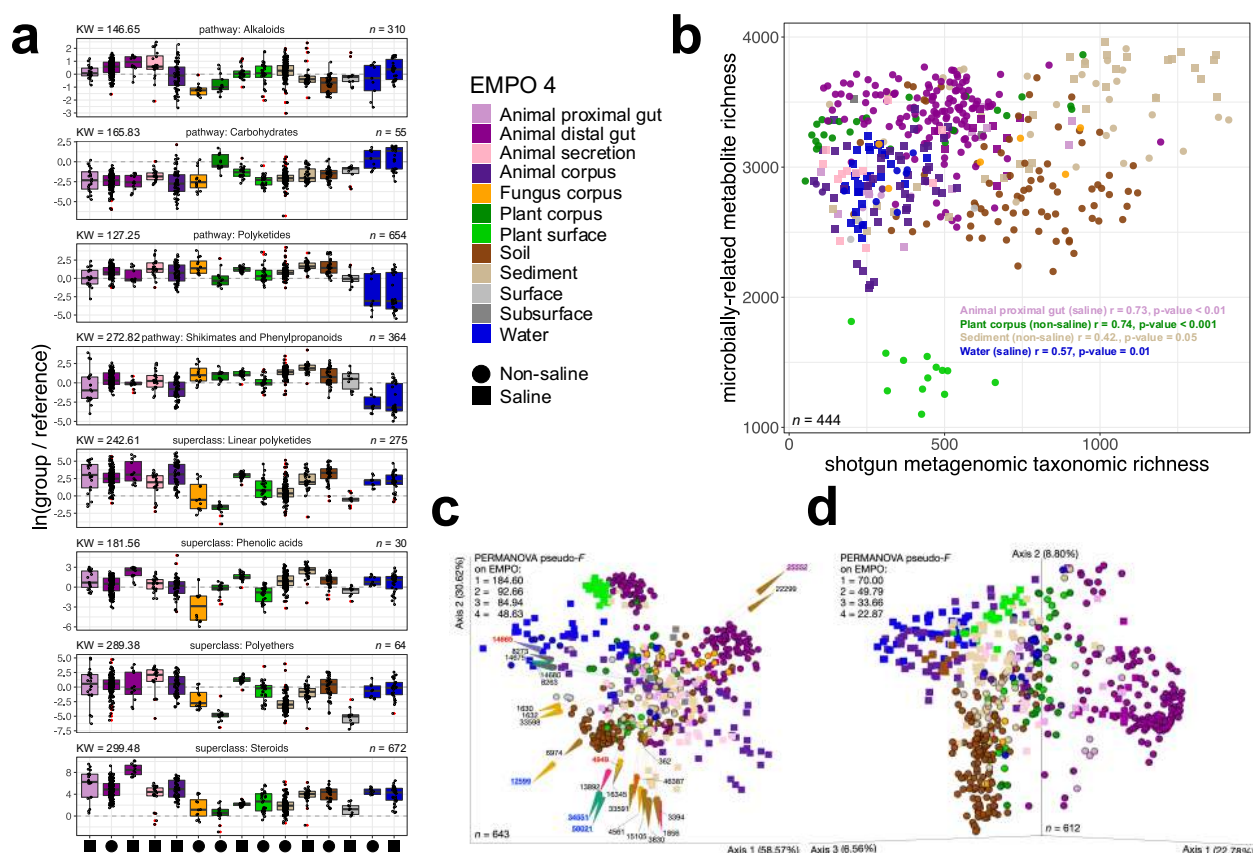
abundances of microbially-related metabolites in aquatic samples had a strong correlation with microbe-metabolite co-occurrences (Figure 3b,c). In addition to highlighting such environments as unique vs. other free-living and host-associated samples, this demonstrates that microbes and metabolites can be classified by- and co-occur among environments.

We further generated additional data for EMP500 samples, including gas chromatography-mass spectrometry (GC-MS) and amplicon sequence (i.e., 16S, 18S, ITS, and full-length rRNA operon) data that also supported a strong relationship between bacterial- and archaeal communities and metabolic profiles (Extended Data Table 1). We anticipate that advances in genome-mining will improve the discovery and classification of BGCs from MAGs and provide additional insight into these findings, and by making these data publicly available in Qiita and GNPS our data will provide an important resource for continued collaborative investigations. In the same manner, the development of novel instrumentation and computational methods for metabolomics will expand the depth of metabolites surveyed in microbiome studies.





# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

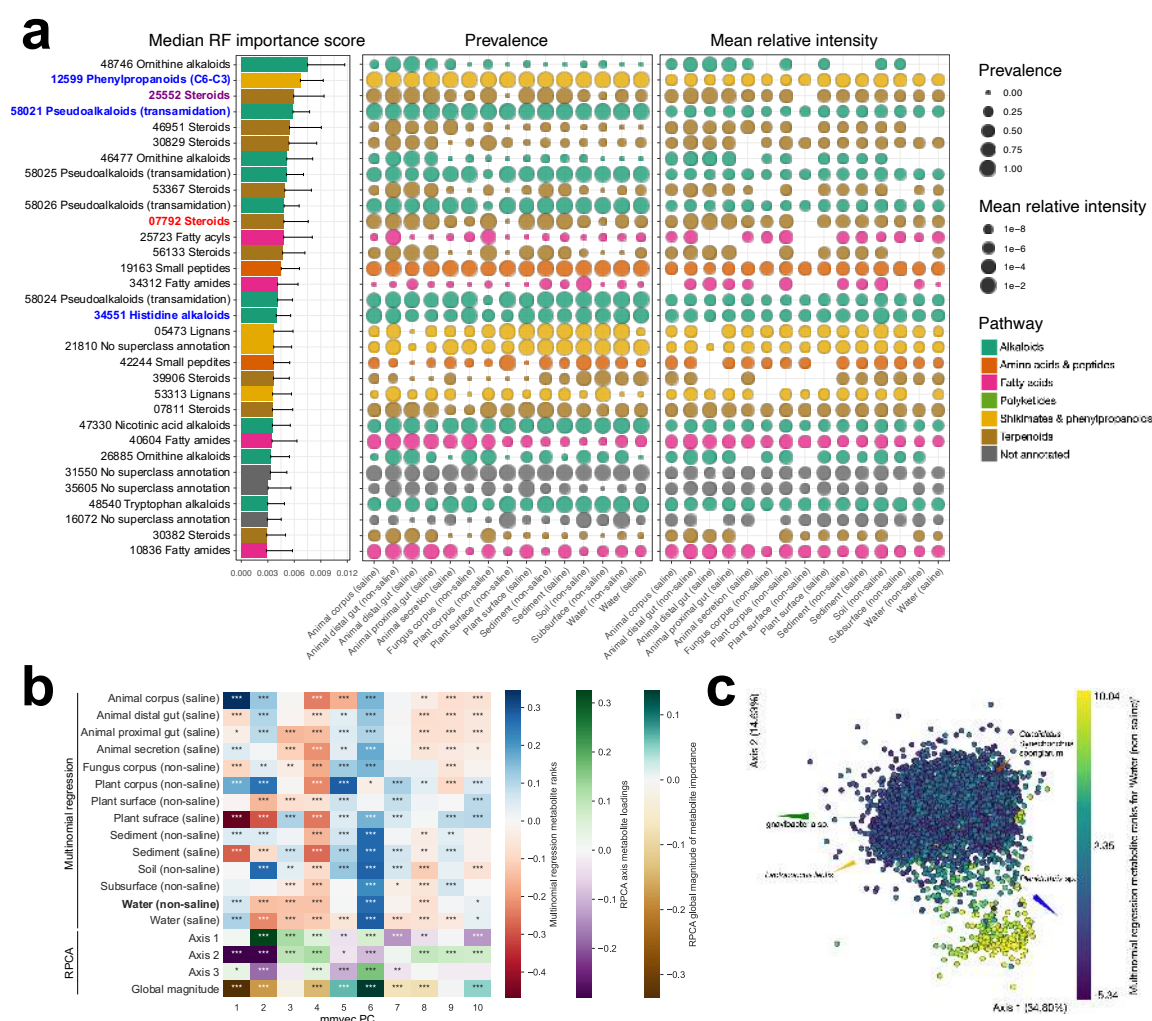


**Fig. 2 | Structural-level associations between microbially-related secondary metabolites and specific environments described using the Earth Microbiome Project Ontology (EMPO version 2, level 4). a,** Differential abundance of molecular features across environments, highlighting four example pathways and four superclasses in separate panels. For each panel, the y-axis represents the natural log-ratio of the intensities of metabolites annotated as the listed ingroup divided by the intensities of metabolites annotated as the reference group (i.e., *Amino Acids and Peptides*,  $n = 615$ , for pathways and *Flavonoids*,  $n = 42$ , for superclasses). The number of metabolites in each ingroup is shown, as well as the chi-squared statistic from a Kruskal-Wallis rank sum test for differences in log-ratios across environments (i.e., each test had  $p\text{-value} < 2.2\text{E-}16$ ). Each test included 606 samples. Outliers from boxplots are colored red to highlight that they are also represented in the overlaid, jittered points. Associations between molecular features and environments were identified using Songbird multinomial regression (model: *composition* = EMPO version 2, level 4; pseudo- $Q^2 = 0.21$ ). Additional information about features is described in Table S2. **b,** Relationship between microbially-related metabolite richness and microbial taxon richness across

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

samples and environments, with significant relationships noted. Correlations with metabolite richness were weaker when using Faith's PD and weighted Faith's PD for quantifying microbial alpha-diversity (Table S5). **c**, Turnover in composition of microbially-related secondary metabolites across environments, visualized using Robust Aitchison Principal Components Analysis (RPCA) showing samples separated based on LC-MS/MS spectra. Shapes represent samples and are colored and shaped by EMPO. Arrows represent metabolites, and are colored by chemical pathways. The direction and magnitude of each arrow corresponds to the strength of the correlation between the relative abundance (i.e., intensity) of the metabolite it represents and the RPCA axes. Samples close to arrow heads have strong, positive associations with respective features, whereas samples at and beyond arrow origins strong, negative associations. The 25 most important metabolites are shown and are described in Table S3. Features annotated in red are those also identified in our multinomial regression analysis as among the top 10 ranked metabolites per environment (Tables S2), those in blue also separated environments in machine-learning analysis (Table S4), and those in purple identified as important in all three analyses. **d**, Turnover in composition of microbial taxa across environments, visualized using Principal Coordinates Analysis (PCoA) of weighted UniFrac distances. Distances are based on counts of microbial genomes from mapping metagenomic reads to the Web of Life database. Note the similarities between panels **c** and **d** with respect to the separation of free-living (e.g., 'Water') and host-associated (e.g., 'Animal distal gut') environments along Axis 1, and a gradient from living hosts (e.g., 'Plant surface'), to detritus (e.g., 'Plant corpus'), to soils and sediments along Axis 2. Results from PERMANOVA for each level of EMPO are shown (i.e., all tests had  $p$ -value = 0.001).

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment



**Fig. 3 | Microbially-related metabolites classify environments with a high accuracy and co-occur with specific microbial taxa. a,** The random-forest importance score, environment-wide prevalence, and mean relative intensity for the top 32 most important microbially-related metabolites contributing to the separation of environments. Metabolites are further described in Table S4. Those in red are those also identified in our multinomial regression analysis as among the top 10 ranked metabolites per environment (Tables S2), those in blue also identified to be strongly associated with RPCA axes (Fig. 2c, Table S3). **b,** Co-occurrence analysis results showing correlation between *mmvec* principal coordinates (PCs) and (i) multinomial regression betas for metabolite abundances across environments, (ii) axes from the RPCA biplot in Fig. 2d corresponding to clustering of samples by environment (i.e., EMPO version 2, level 4) based on metabolite profiles, and (iii) a vector representing the global magnitude of features from the RPCA biplot in Fig. 2d. Values are Spearman correlation coefficients. Asterisks indicate significant

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

correlations ( $*p < 0.05$ ,  $**p < 0.01$ ,  $***p < 0.001$ ). **c**, The relationship between multinomial regression betas for metabolites with respect to ‘Water (non-saline)’ and the first three *mmvec* PCs shown as a multi-omics biplot of metabolite-microbe co-occurrences learned from microbial profiles. Points represent metabolites separated by their co-occurrences with microbial taxa. Metabolites are colored based on multinomial regression betas for their abundances with respect to ‘Water (non-saline)’. Vectors represent specific microbial taxa strongly associated with ordination axes. The model predicting metabolite-microbe co-occurrences was more accurate than one representing a random baseline, with a pseudo- $Q^2$  value of 0.18, indicating much reduced error during cross-validation.

## References.

1. Davies, J. Specialized microbial metabolites: functions and origins. *J. Antibiot.* **66**, 361–364 (2013).
2. Pham, J. V. *et al.* A Review of the Microbial Production of Bioactive Natural Products and Biologics. *Front. Microbiol.* **10**, 1404 (2019).
3. Nayfach, S. *et al.* A genomic catalog of Earth’s microbiomes. *Nat. Biotechnol.* (2020)  
doi:10.1038/s41587-020-0718-6.
4. Ziemert, N., Alanjary, M. & Weber, T. The evolution of genome mining in microbes - a review. *Nat. Prod. Rep.* **33**, 988–1005 (2016).
5. Libis, V. *et al.* Uncovering the biosynthetic potential of rare metagenomic DNA using co-occurrence network analysis of targeted sequences. *Nat. Commun.* **10**, 3848 (2019).
6. Hugerth, L. W. *et al.* Metagenome-assembled genomes uncover a global brackish microbiome. *Genome Biol.* **16**, 279 (2015).
7. Tully, B. J., Graham, E. D. & Heidelberg, J. F. The reconstruction of 2,631 draft metagenome-assembled genomes from the global oceans. *Sci Data* **5**, 170203 (2018).
8. Stewart, R. D. *et al.* Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery. *Nat. Biotechnol.* **37**, 953–961 (2019).
9. Van Goethem, M. W. *et al.* Long-read metagenomics of soil communities reveals phylum-specific secondary metabolite dynamics. *Cold Spring Harbor Laboratory* 2021.01.23.426502 (2021)  
doi:10.1101/2021.01.23.426502.
10. Amos, G. C. A. *et al.* Comparative transcriptomics as a guide to natural product discovery and biosynthetic gene cluster functionality. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E11121–E11130 (2017).
11. Vogel, C. & Marcotte, E. M. Insights into the regulation of protein abundance from proteomic and transcriptomic analyses. *Nat. Rev. Genet.* **13**, 227–232 (2012).
12. Aksenov, A. A., da Silva, R., Knight, R., Lopes, N. P. & Dorrestein, P. C. Global chemical analysis of biology by mass spectrometry. *Nature Reviews Chemistry* **1**, 0054 (2017).

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

13. Dührkop, K. *et al.* Systematic classification of unknown metabolites using high-resolution fragmentation mass spectra. *Nat. Biotechnol.* (2020) doi:10.1038/s41587-020-0740-8.
14. Mohimani, H. *et al.* Dereplication of microbial metabolites through database search of mass spectra. *Nat. Commun.* **9**, 4035 (2018).
15. Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).
16. van Santen, J. A. *et al.* The Natural Products Atlas: An Open Access Knowledge Base for Microbial Natural Products Discovery. *ACS Cent Sci* **5**, 1824–1833 (2019).
17. Kautsar, S. A. *et al.* MIBiG 2.0: a repository for biosynthetic gene clusters of known function. *Nucleic Acids Res.* **48**, D454–D458 (2020).
18. Nothias, L.-F. *et al.* Feature-based molecular networking in the GNPS analysis environment. *Nat. Methods* **17**, 905–908 (2020).

## Acknowledgements.

We thank Lindsay Goldasich and Julia Toronczak for assistance with sample processing for sequencing; Marcus Fedarko, Rachel Diner, Joshua Ladau, Elisha Wood-Charlson, Stephen Nayfach, Daniel Udvary, and Emiley Eloie-Fadrosch for reviewing an early version of the manuscript. This work was supported in part by the United States (US) National Institute of Health (NIH) (awards 1RF1-AG058942-01, 1DP1AT010885, R01HL140976, R01DK102932, R01HL134887, U19AG063744, U01 AI124316), US Department of Agriculture – National Institute of Food and Agriculture (USDA-NIFA) (award 2019-67013-29137), the US National Science Foundation (NSF) - Center for Aerosol Impacts on Chemistry of the Environment, Crohn’s & Colitis Foundation Award (CCFA) (award 675191), Semiconductor Research Corporation and Defense Advanced Research Projects Agency (SRC/DARPA) (award GI18518), Department of Defense (award W81XWH-17-1-0589), the Office of Naval Research (ONR) (award N00014-15-1-2809), the Emerald Foundation (award 3022), IBM Research AI through the AI



## Shaffer et al. Metabolite-microbe profiles are shaped by the environment

Horizons Network, and the Center for Microbiome Innovation. Metabolomics analyses at Pacific Northwest National Laboratory (PNNL) were supported by the Laboratory Directed Research and Development program via the Microbiomes in Transition Initiative and performed in the Environmental Molecular Sciences Laboratory, a national scientific user facility sponsored by the U.S. Office of Biological and Environmental Research and located at PNNL. PNNL is a multiprogram national laboratory operated by Battelle for the Department of Energy (DOE) under contract DE-AC05-76RLO 1830. J.P.S. was supported by NIH/NIGMS IRACDA K12 GM068524. L.F.N. was supported by the NIH (award R01-GM107550). A.D.B. was supported by the Danish Council for Independent Research (DFF) (award 9058-00025B). W.B. was supported by the Research Foundation – Flanders (12W0418N). K.D. and S.B. were supported by Deutsche Forschungsgemeinschaft (BO 1910/20 and 1910/23). J.S. was supported by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research, Environmental System Science (ESS) Program, and his contribution originates from the River Corridor Scientific Focus Area (SFA) project at PNNL. P.C.D. was supported by the Gordon and Betty Moore Foundation (award GBMF7622) and the NIH (award R01-GM107550). We thank Eppendorf, Illumina, and Integrated DNA Technologies for in-kind support at various phases of the project.

### **Author contributions.**

J.A.G., J.K.J., and R.K. conceived the idea for the project. P.C.D., and R.K. designed the multi-omics component of the project and provided project oversight. J.P.S. managed the project, performed preliminary data exploration, coordinated data analysis, analyzed data, and provided data interpretation. L.F.N. coordinated and performed LC-MS/MS sample processing, coordinated and performed annotation of LC-MS/MS, analyzed LC-MS/MS data, and provided LC-MS/MS data interpretation. L.R.T. designed the multi-omics component of the project, solicited sample collection, curated sample metadata, processed samples, performed preliminary data exploration, and provided project oversight. J.G.S. designed the multi-omics component, managed the project, developed protocols and tools, coordinated and performed sequencing, and performed preliminary exploration of sequence data. R.A.S. developed



# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

protocols and coordinated and performed sequencing. S.P.C. and T.O.M. coordinated and performed GC-MS sample processing and provided interpretation of GC-MS data. A.D.B. conceived the idea for the paper, performed preliminary data exploration, analyzed data, and provided data interpretation. S.H. performed machine-learning analyses. F.L. performed co-occurrence analysis, multinomial regression analyses, and correlations with co-occurrence data. H.L.L. performed multinomial regression analyses. Q.Z. developed tools and provided interpretation of shotgun metagenomics data. C.Mart. and J.T.M. provided oversight and interpretation of RPCA, multinomial regression, and co-occurrence analyses. S.K. performed preliminary exploration of shotgun metagenomics data. K.D., S.B, and H.W.K. annotated LC-MS/MS data. A.A.A. processed GC-MS data. W.B. provided oversight for machine-learning analyses. C.Maro. processed samples for sequencing. Y.V.B. performed preliminary data exploration and provided oversight for machine-learning analysis. A.T. and D.P. performed preliminary data exploration. L.P., A.P.C., N.H., and K.L.B. performed preliminary exploration of shotgun metagenomic data and performed machine learning analyses. P.D. performed preliminary exploration of shotgun metagenomics data. A.G. developed tools, provided interpretation of shotgun metagenomics data, and analyzed shotgun metagenomics data. G.H. coordinated short-read amplicon and shotgun metagenomics sequencing. M.M.B. and K.S. performed short-read amplicon and shotgun metagenomics sequencing. T.S. assisted with DNA extraction. D.M. coordinated long-read amplicon sequencing, analyzed shotgun metagenomics data, and provided interpretation of the data. S.M.K. and M.A. coordinated and performed long-read amplicon sequencing and long-read sequence data analysis. J.J.M. collected samples, coordinated field logistics, developed protocols, and performed short-read amplicon and shotgun metagenomics sequencing. S.S. collected samples, coordinated field logistics, and provided interpretation of the data. G.L.A. curated sample metadata and organized sequence data. J.D. processed sequence data. A.D.S. provided project oversight and data interpretation. T.T., A.S., and J.S. collected samples, coordinated field logistics, and provided interpretation of the data. J.P.S. and L.F.N. wrote the manuscript, with contributions from all authors.

Shaffer et al. Metabolite-microbe profiles are shaped by the environment

# **Earth Microbiome Project 500 (EMP500) Consortium.**

Jennifer L. Bowen<sup>1</sup>, Max Chavarría-Vargas<sup>2</sup>, Don A. Cowan<sup>3</sup>, Jaime Huerta-Cepas<sup>4</sup>, Paul Jensen<sup>5</sup>,  
Lingjing Jiang<sup>6</sup>, Anton Lavrinienko<sup>7</sup>, Thulani P. Makhalanyane<sup>3</sup>, Tapio Mappes<sup>7</sup>, Ezequiel M.  
Marzinelli<sup>8</sup>, Sou Miyake<sup>9</sup>, Timothy A. Mousseau<sup>7</sup>, Catalina Murillo-Cruz<sup>2</sup>, Adrián A. Pinto-Tomás<sup>2</sup>,  
Dorota L. Porazinska<sup>10</sup>, Jean-Baptiste Ramond<sup>3,11</sup>, Taniya RoyChowdhury<sup>12,13</sup>, Henning Seedorf<sup>9,14</sup>, J.  
Reuben Shipway<sup>15,16</sup>, Frank J. Stewart<sup>17</sup>, Jana M. U'Ren<sup>18</sup>, Phillip C. Watts<sup>7</sup>, Nicole S. Webster<sup>19,20</sup>, Jesse  
R. Zaneveld<sup>21</sup>, Shan Zhang<sup>22</sup>

<sup>1</sup>Northeastern University, Boston, Massachusetts, USA. <sup>2</sup>University of Costa Rica, San José, Costa Rica.  
<sup>3</sup>University of Pretoria, Pretoria, South Africa. <sup>4</sup>Universidad Politécnica de Madrid, Instituto Nacional de  
Investigación y Tecnología Agraria y Alimentaria, Madrid, Spain. <sup>5</sup>University of California San Diego,  
La Jolla, California, USA. <sup>6</sup>Janssen Research & Development, San Diego, California, USA. <sup>7</sup>University  
of Jyväskylä, Jyväskylä, Finland. <sup>8</sup>The University of Sydney, Sydney, Australia. <sup>9</sup>Temasek Life Sciences  
Laboratory, Singapore, Singapore. <sup>10</sup>University of Florida, Gainesville, Florida, USA. <sup>11</sup>Pontificia  
Universidad Católica de Chile, Santiago, Chile. <sup>12</sup>Pacific Northwest National Laboratory, Richland,  
Washington, USA. <sup>13</sup>University of Maryland, College Park, Maryland, USA. <sup>14</sup>National University of  
Singapore, Singapore, Singapore. <sup>15</sup>University of Portsmouth, Portsmouth, United Kingdom. <sup>16</sup>University  
of Massachusetts Amherst, Amherst, Massachusetts, USA. <sup>17</sup>Montana State University, Bozeman,  
Montana, USA. <sup>18</sup>University of Arizona, Tucson, Arizona, USA. <sup>19</sup>Australian Institute of Marine Science,  
Townsville, Qld, Australia. <sup>20</sup>University of Queensland, Brisbane, Australia. <sup>21</sup>University of Washington  
Bothell, Bothell, Washington, USA. <sup>22</sup>University of New South Wales, Sydney, Australia.

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment

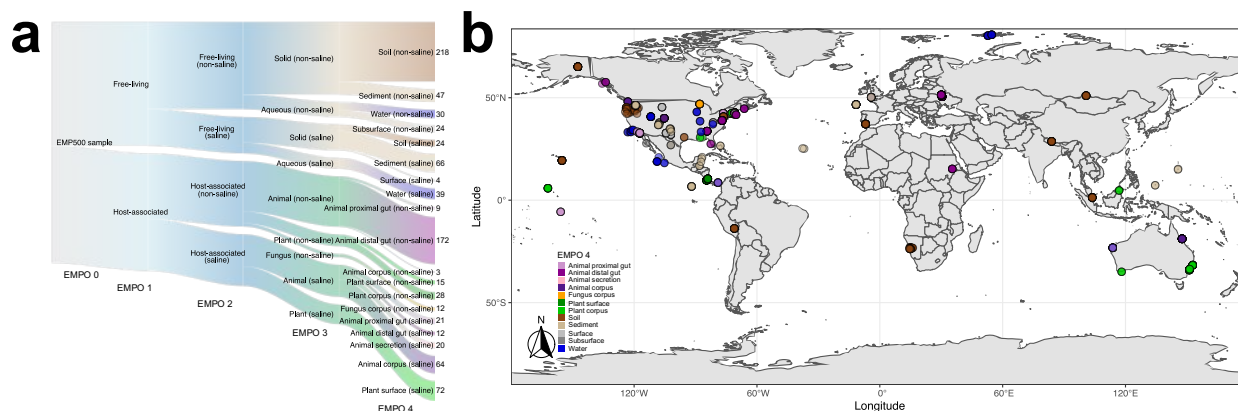
## Extended Data.

### Extended Data Table 1 | Mantel test results comparing data layers generated for the EMP500 samples.

Note the strong relationships between the metabolomics data (i.e., LC-MS/MS and GC-MS) and the sequence data from Bacteria and Archaea (i.e., shotgun metagenomics, 16S, and full-length rRNA operon) as compared to between metabolomics data and sequence data from eukaryotes (i.e., 18S and ITS), as well as the strong relationships between difference sequence data from Bacteria and Archaea (rho > 0.2 in bolded font; > 0.4 in bolded, italics; > 0.5 additionally underlined).

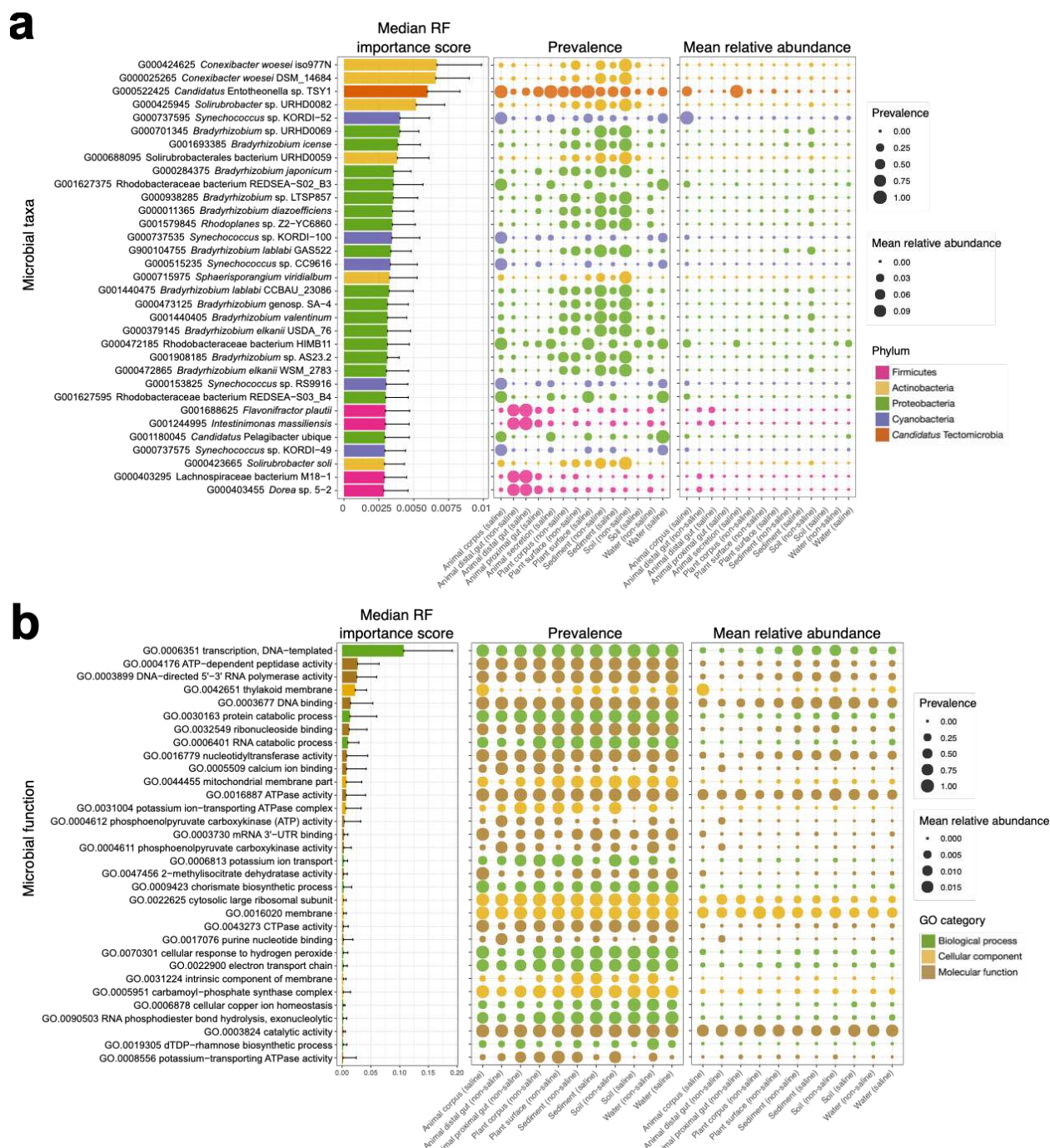
Dataset 1	Dataset 2	<i>n</i>	Spearman rho	p-value	
LC-MS/MS	GC-MS	401	0.13	0.001	
	Metagenomics (taxa)	454	<b>0.43</b>	0.001	
	Metagenomics (function)	378	<b>0.24</b>	0.001	
	16S	477	<b>0.27</b>	0.001	
	18S	356	0.06	0.034	
	ITS	359	0.08	0.001	
	full-length rRNA operon	181	<b>0.34</b>	0.001	
GC-MS	Metagenomics (taxa)	331	0.07	0.002	
	Metagenomics (function)	279	0.05	0.029	
	16S	349	<b>0.22</b>	0.001	
	18S	284	-0.01	0.641	
	ITS	258	0.02	0.404	
	full-length rRNA operon	168	0.11	0.001	
	Metagenomics (taxa)	Metagenomics (function)	508	<b>0.24</b>	0.001
	16S	538	<b>0.51</b>	0.001	
	18S	378	-0.04	0.159	
	ITS	408	0.05	0.023	
	full-length rRNA operon	235	<b>0.48</b>	0.001	
	Metagenomics (function)	16S	449	0.11	0.001
	18S	305	0.12	0.001	
	ITS	337	0.14	0.001	
	full-length rRNA operon	219	<b>0.21</b>	0.001	
	16S	18S	430	0.04	0.117
	ITS	447	0.09	0.001	
	full-length rRNA operon	215	<b>0.51</b>	0.001	
	18S	ITS	387	<b>0.20</b>	0.001
	full-length rRNA operon	174	0.06	0.06	
ITS	full-length rRNA operon	166	0.09	0.001	

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment



**Extended Data Fig. 1 | a**, Distribution of samples among the Earth Microbiome Project Ontology (EMPO version 2) categories. **b**, Geographic distribution of samples with points colored by EMPO (version 2, level 4). Extensive information about each sample set contributed is described in Table S1.

# Shaffer et al. Metabolite-microbe profiles are shaped by the environment



**Extended Data Fig. 2 | Machine-learning analysis of shotgun metagenomics data highlighting the**

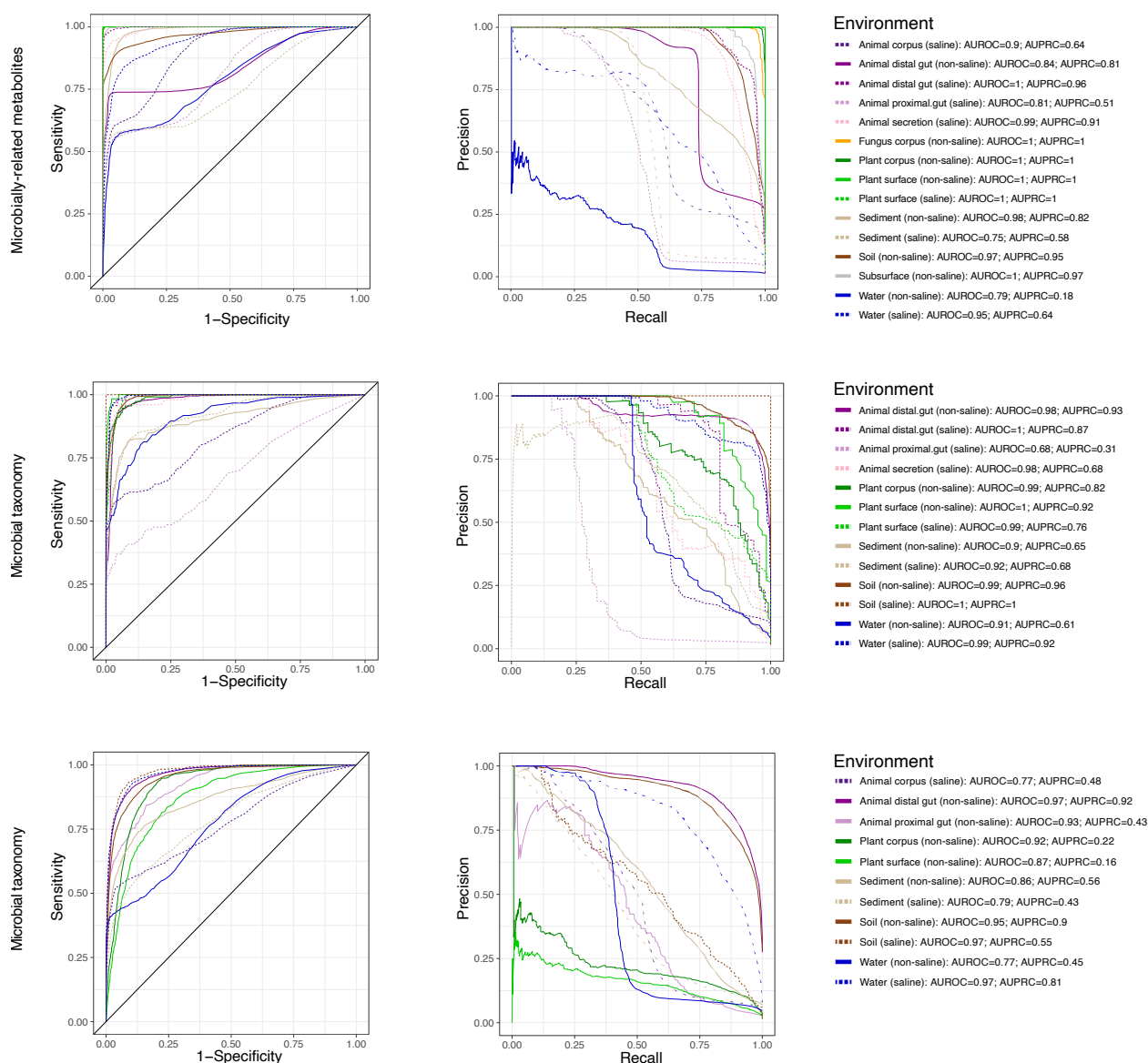
**most important microbial taxa and functions.** **a**, The random-forest (RF) importance score, environment-wide prevalence, and mean relative abundance for the top 32 most important microbial taxa contributing to the separation of environments. **b**, The RF importance score, environment-wide prevalence, and mean relative abundance for the top 32 most important microbial functions (i.e., GO

## Shaffer et al. Metabolite-microbe profiles are shaped by the environment

Terms) contributing to the separation of environments. For both analyses, environments are described by the Earth Microbiome Project Ontology (EMPO version 2, level 4). All samples were analyzed using 20-time repeated group five-fold cross validation, where ‘group’ indicates the study identifier that a sample belongs to. In this approach, we split the data by ‘group’ such that all samples from the same study identifier can be assigned to either train or test data sets.



# Shaffer et al. Metabolite-microbe profiles are shaped by the environment



## Extended Data Fig. 3 | Machine-learning analysis of LC-MS/MS metabolomics and shotgun

## metagenomics data highlighting per-environment classification accuracy. **a**, Receiver operating

characteristic curves (and AUROC) and precision-recall curves (and AUPRC) illustrating classification accuracy of the random forest model across all environments based on microbial taxonomic (i.e., OGU) profiles. **b**, Receiver operating characteristic curves (and AUROC) and precision-recall curves (and

AUPRC) illustrating classification accuracy of the random forest model across all environments based on

microbial functional (i.e., GO terms) profiles. For both analyses, environments are described by the Earth Microbiome Project Ontology (EMPO version 2, level 4). All samples were analyzed using 20-time

## Shaffer et al. Metabolite-microbe profiles are shaped by the environment

repeated group five-fold cross validation, where ‘group’ indicates the study identifier that a sample belongs to. In this approach, we split the data by ‘group’ such that all samples from the same study identifier can be assigned to either train or test data sets.