

1 **Title:** Reference data based insights expand 2 understanding of human metabolomes

3 Authors: Julia M. Gauglitz^{1,2}, Wout Bittremieux^{1,2,3}, Candace L. Williams⁴, Kelly C. Weldon^{1,2,5},
4 Morgan Panitchpakdi^{1,2}, Francesca Di Ottavio¹, Christine M. Aceves^{1,2}, Elizabeth Brown^{2,6}, Nicole
5 C. Sikora^{1,2}, Alan K. Jarmusch^{1,2}, Cameron Martino^{5,7,8}, Anupriya Tripathi^{2,6,7}, Erfan Sayyari⁵,
6 Justin P. Shaffer⁷, Roxana Coras⁹, Fernando Vargas^{1,2,6}, Lindsay DeRight Goldasich⁷, Tara
7 Schwartz⁷, MacKenzie Bryant⁷, Gregory Humphrey⁷, Abigail J. Johnson¹⁰, Katharina Spengler¹,
8 Pedro Belda-Ferre⁷, Edgar Diaz⁷, Daniel McDonald⁷, Qiyun Zhu⁷, Dominic S. Nguyen⁷, Emmanuel
9 O. Elijah^{1,2}, Mingxun Wang^{1,2}, Clarisse Marotz⁷, Kate E. Sprecher^{11,12}, Daniela Vargas Robles¹³,
10 Dana Withrow¹¹, Gail Ackermann⁷, Lourdes Herrera¹⁴, Barry J. Bradford¹⁵, Lucas Maciel Mauriz
11 Marques¹⁶, Juliano Geraldo Amaral¹⁷, Rodrigo Moreira Silva¹⁸, Flávio Protaso Veras¹⁶, Thiago
12 Mattar Cunha¹⁶, Rene Donizeti Ribeiro Oliveira¹⁹, Paulo Louzada-Junior¹⁹, Robert H. Mills^{1,2,7,20},
13 Douglas Galasko²¹, Parambir S. Dulai²², Curt Wittenberg²³, David J. Gonzalez^{1,2,5,20}, Robert
14 Terkeltaub²¹, Megan M. Doty^{7,24}, Jae H. Kim²⁵, Kyung E. Rhee⁷, Julia Beauchamp-Walters²⁶,
15 Kenneth P. Wright Jr¹¹, Maria Gloria Dominguez-Bello²⁷, Mark Manary²⁸, Michelli F. Oliveira²⁹,
16 Brigid S. Boland²¹, Norberto Peporine Lopes¹⁸, Monica Guma²¹, Austin D. Swafford⁵, Rachel J.
17 Dutton⁶, Rob Knight^{5,7,30,31,*}, Pieter C. Dorrestein^{1,2,5,7,*}

18
19

20 **Author Affiliations:**

21 1 Collaborative Mass Spectrometry Innovation Center; University of California San Diego; La
22 Jolla, CA 92093; USA

23 2 Skaggs School of Pharmacy and Pharmaceutical Sciences; University of California San Diego;
24 La Jolla CA 92093; USA

25 3 Department of Computer Science; University of Antwerp; 2020 Antwerpen; Belgium

26 4 Reproductive Sciences, Institute for Conservation Research; San Diego Zoo Global; Escondido,
27 CA 92027; USA

28 5 Center for Microbiome Innovation, Joan and Irwin Jacobs School of Engineering; University of
29 California San Diego; La Jolla, CA 92093; USA

30 6 Division of Biological Sciences; University of California San Diego; La Jolla, CA 92093; USA

31 7 Department of Pediatrics, School of Medicine; University of California San Diego; La Jolla, CA
32 92093; USA

33 8 Bioinformatics and Systems Biology Program; University of California San Diego; La Jolla, CA
34 92093; USA

35 9 Division of Rheumatology, Allergy & Immunology, Department of Medicine; University of
36 California San Diego; La Jolla, CA 92093; USA

37 10 BioTechnology Institute; University of Minnesota; Saint Paul, MN 55108; USA

38 11 Department of Integrative Physiology; University of Colorado Boulder; Boulder, CO 80309;
39 USA

40 12 Department of Population Health Sciences; University of Wisconsin-Madison; Madison, WI
41 53726; USA

- 42 13 Servicio Autónomo Centro Amazónico de Investigación y Control de Enfermedades Tropicales
43 Simón Bolívar; Puerto Ayacucho 7101, Amazonas; Venezuela
- 44 14 Department of Pediatrics; Wake Forest School of Medicine; Winston-Salem, NC 27101; USA
- 45 15 Department of Animal Science; Michigan State University; East Lansing, MI 48824; USA
- 46 16 Department of Pharmacology, School of Medicine of Ribeirão Preto, Center of Research in
47 Inflammatory Diseases; University of São Paulo; Ribeirão Preto, CEP 14049-900 - SP; Brazil
- 48 17 Multidisciplinary Health Institute; Federal University of Bahia; 45029094, Vitória da Conquista
49 - BA; Brazil
- 50 18 NPPNS, Department of Biomolecular Sciences, School of Pharmaceutical Sciences of
51 Ribeirão Preto; University of São Paulo; Ribeirão Preto. CEP 14040-903 - SP; Brazil
- 52 19 Department of Internal Medicine, School of Medicine of Ribeirão Preto, Center of Research in
53 Inflammatory Diseases; University of São Paulo; Ribeirão Preto, CEP 14049-900 - SP; Brazil
- 54 20 Department of Pharmacology; University of California San Diego; La Jolla, CA 92093; USA
- 55 21 Department of Neurosciences; University of California San Diego; La Jolla, CA 92093; USA
- 56 22 Division of Gastroenterology, Department of Medicine; University of California San Diego; La
57 Jolla, CA 92093; USA
- 58 23 Department of Molecular Medicine; The Scripps Research Institute; La Jolla, CA 92037; USA
- 59 24 Division of Neonatology, Department of Pediatrics, Kapi'olani Medical Center for Women and
60 Children; John A. Burns School of Medicine; Honolulu, Hawaii 96813; USA
- 61 25 Division of Neonatology, Perinatal Institute, Department of Pediatrics, Cincinnati Children's
62 Hospital Medical Center; University of Cincinnati College of Medicine; Cincinnati, Ohio 45229;
63 USA
- 64 26 Division of Pediatric Hospital Medicine, School of Medicine; University of California San Diego;
65 La Jolla, CA 92093; USA
- 66 27 Department of Biochemistry and Microbiology, School of Environmental and Biological
67 Sciences; Rutgers, The State University of New Jersey; New Brunswick, NJ 08901; USA
- 68 28 Department of Pediatrics; Washington University; St. Louis, MO 63110; USA
- 69 29 Department of Medicine; University of California San Diego; La Jolla, CA 92093; USA
- 70 30 Department of Computer Science and Engineering; University of California San Diego; La
71 Jolla, CA 92093; USA
- 72 31 Department of Bioengineering; University of California San Diego; La Jolla, CA 92093; USA
73

74 Summary

75 The human metabolome has remained largely unknown, with most studies annotating ~10% of
76 features. In nucleic acid sequencing, annotating transcripts by source has proven essential for
77 understanding gene function. Here we generalize this concept to stool, plasma, urine and other
78 human metabolomes, discovering that food-based annotations increase the interpreted fraction
79 of molecular features 7-fold, providing a general framework for expanding the interpretability of
80 human metabolomic “dark matter.”
81
82

83 Introduction

84 In 2016, typical MS/MS-based untargeted metabolomics studies annotated only ~2% of
85 molecules based on matches against spectral libraries, leaving the rest of the sample as
86 metabolomic “dark matter.” The capture of community knowledge, accumulating public reference
87 MS/MS spectra over the past four years, has increased this baseline ~2.5-fold within the global
88 natural product social molecular networking (GNPS) infrastructure (Wang et al., 2016). This
89 growth has been even more dramatic for data from commonly-studied specimen types such as
90 human stool and plasma: 10.1 +/- 4.4% of MS/MS features now match to a reference MS/MS
91 spectrum [1% FDR (Scheubert et al., 2017), n = 30, average number of unique MS/MS spectra is
92 12,889/dataset]. However, despite these advances, the vast majority of detectable spectra lack
93 any annotation.

94 This situation for MS/MS spectra is in sharp contrast to the interpretability of
95 uncharacterized portions of the human genome. For example, reference data sets for gene
96 expression, such as expressed sequence tags (an early form of RNASeq), enable the sequencing
97 of “dark matter,” as opposed to monitoring the expression of a single curated gene. Such methods
98 have significantly improved interpretation by annotating genes not directly by function, but rather
99 by source (developmental stage, tissue location, organism-level, phenotype, etc.) (Bono, 2020;
100 Ono et al., 2017). Interpretation based on source has been very important for metagenomics and
101 metatranscriptomics, increasing our understanding of the structure and function of complex
102 communities by leveraging matches between genes or transcripts of known and unknown origin
103 via publicly available databases.

104 Annotation of chemicals, based on their source within publicly available complex reference
105 samples that use controlled metadata vocabularies, has not been applied to metabolomics for
106 several reasons. First, standards for annotation of molecules that are used to create spectral
107 libraries have been based on availability of individual pure, typically commercially available,
108 standards, and structural considerations such as presence of specific moieties. Many molecules
109 are observed as multiple different ion forms, such as adducts, in-source fragments, and
110 multimers. Current spectral libraries do not contain all possible ion forms of those molecules, and
111 typically only the protonated form (Schmid et al., 2020; Vinaixa et al., 2016), because reference
112 standards that run in a highly purified state that biases towards detection extraction of only specific
113 data on specific ion forms. These forms are often different from the ions associated with the same
114 molecule present in an extract from a biological matrix (e.g. proton vs sodium or even multiple
115 sodium and potassium adducts), which then cannot be matched because the relevant spectra are
116 not in the database. Second, on average, 5–10% of untargeted metabolomics data can be
117 annotated from spectral libraries: the remaining 90+% are unassignable “dark matter” in
118 metabolomics, especially when obtained from complex matrices such as human samples. Third,
119 large databases of untargeted metabolomics data with consistently annotated provenance with
120 controlled vocabularies have been neither available nor possible to effectively reuse. We recently
121 addressed this latter problem via GNPS (Wang et al., 2016), ReDU (Jarmusch et al., 2019),
122 importing data from MetaboLights into GNPS (Haug et al., 2020), with ReDU-compatible
123 metadata conversion. Finally, the availability of robust scalable analysis infrastructures and
124 algorithms, such as molecular networking, that enable the functional equivalent of reporting of

125 expressed sequence tag/RNASeq analysis, have only recently been introduced for mass
126 spectrometry (Wang et al., 2016; Watrous et al., 2012).

127 To improve interpretation of otherwise unannotated data from untargeted mass
128 spectrometry experiments, we leverage entire reference data sets with curated ontologies to
129 complement existing spectral libraries of individual molecules. Due to lack of a better term we
130 refer to this approach as interpretive metabolomics in this manuscript, and demonstrate its
131 potential by leveraging the Global FoodOmics MS/MS spectral database, which we have made
132 publicly available on MassIVE. This food reference data set will be key for enabling future insights
133 into human health given the importance of diet and the urgent need to develop additional methods
134 for empirical nutrient and diet assessments to understand acute and chronic human disease
135 (Barabási et al., 2020). We demonstrate that interpretive metabolomics can address these types
136 of knowledge gaps by showing that it not only massively expands the fraction of the data that can
137 be interpreted, but that these new insights can lead to an improved understanding of the diets
138 consumed upon co-analysis of human and food/beverage mass spectral data.

139 Results/Discussion

140 We conjectured that a major source of chemicals detected by metabolomics in human samples
141 originates from foods and beverages. We created “Global FoodOmics”
142 (<http://www.globalfoodomics.org>) in 2017, which now contains 3,579 food and beverage samples
143 contributed by the community, as outlined in the methods, following in the footsteps of the
144 American Gut and the Earth Microbiome Projects (McDonald et al., 2018; Thompson et al., 2017).
145 The majority of samples were photographed, and a subset were subjected to 16S rRNA profiling
146 (1,511 samples) to characterize the microbial composition, as well as providing information about
147 mitochondria and chloroplast sequences matched by the same primers. Foods were manually
148 classified according to the Earth Microbiome Project Ontology, the USDA Food Composition
149 Database and a modification of the Food and Nutrient Database for Dietary Studies (Johnson et
150 al., 2019; Thompson et al., 2017) (<https://ndb.nal.usda.gov/>) to allow cross-study compatibility. In
151 total, we report 157 metadata categories that further include a six-level food ontology, as well as
152 fermentation or organic status, land or aquatic origin, country of origin, etc. (**Table S1**). Foods
153 and beverages in Global FoodOmics consist of a range of items, from simple ingredients to
154 prepared meals, as well as animal feed.

155 A key benefit of interpretive metabolomics is that we consider all different ion forms
156 encountered while collecting the Global FoodOmics dataset. The millions of MS/MS spectra in
157 Global FoodOmics inherently include MS/MS spectra of different ion forms of both known and
158 unknown molecules, and can, therefore, be matched in human biospecimens via direct matching
159 of the MS/MS spectra or by more sophisticated approaches. The similar complexity of the
160 reference and experimental data includes many chemicals that may have uncharacterized
161 behavior, such as unexpected adducts or even multimers made up of different molecules. For the
162 MS/MS spectra that do have annotations, it is possible to leverage GNPS tags to test whether the
163 spectral matches make sense in the context of Global FoodOmics.

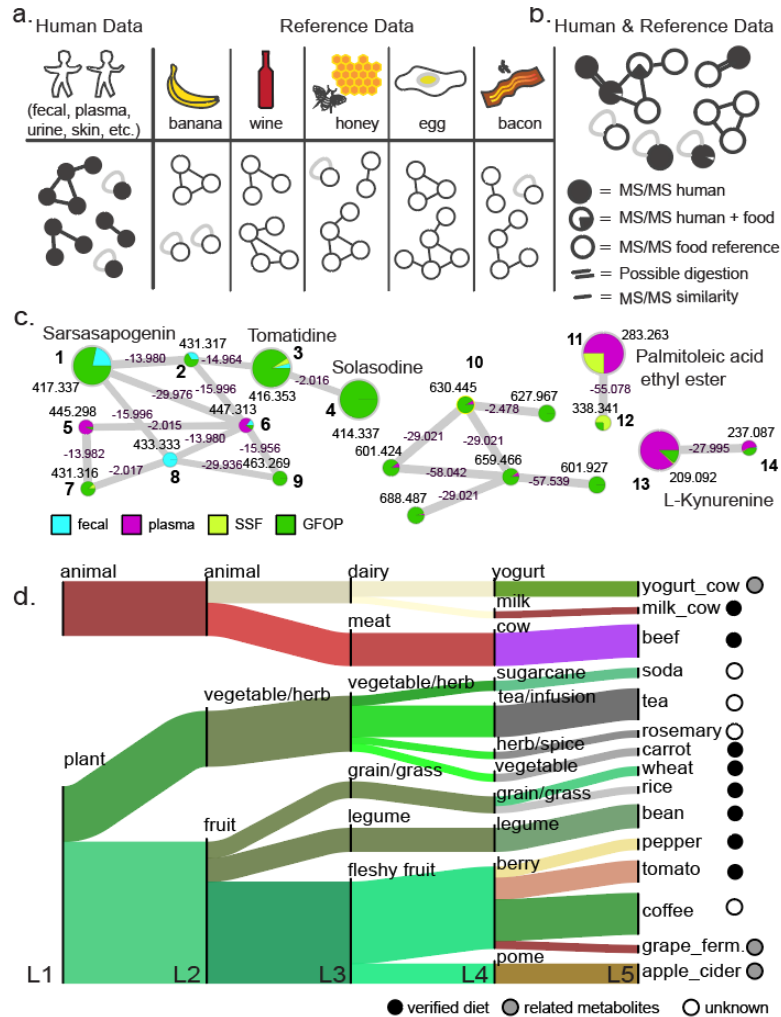
164 Within the GNPS environment, the community can also add tags to each reference
165 spectrum in the spectral library using a controlled vocabulary, including multiple per structure. An
166 InChIKey was included for 4586 of 5455 spectral matches against the reference libraries (~5%

167 annotation rate at 1% FDR), which yielded 1492 unique structures upon consideration of planar
168 structures. There were 415/1492 structures that had lifestyle tags and “food consumption” is the
169 most frequently reported with 357 entries (86%) (**Figure S1a**) (Bouslimani et al., 2016). Brief
170 descriptive tags provide more detail about the annotation itself, and 1131/1492 structures were
171 annotated with such tags. The most common descriptive tags were in order: “natural product”
172 (790/1131), “food” (576/1131), “human”, “plant”, “natural product_plant”, “plant_angiospermae”,
173 and “drug” (**Figure S1b**). Some of these associations with the category “human” may also be of
174 food origin, such as arachidonoyl carnitine, which is currently only tagged as “human,” but may
175 have a variety of animal-product based food sources. Similarly, the tag “drug” includes
176 annotations such as the antimicrobial agent monensin, which is not tagged as a food molecule,
177 but is consumed with animal products from animals raised using monensin as a growth promoter.
178 Thus the Global FoodOmics reference data capture not only inherently food-derived molecules,
179 but also food-sourced exogenous compounds such as preservatives, growth enhancing
180 substances, antimicrobials, pesticides, and packaging materials. However, because the
181 annotation rate remains low, most of the data remains unused despite the informative tags.

182 In addition to annotating molecules based on matches to library spectra, spectral matches
183 to the food reference data can be obtained and visualized using MS/MS based molecular
184 networking. When applying this method to both foods and biospecimens in an experimental sleep
185 restriction and circadian misalignment study we observed connectivity of nodes within molecular
186 families representing MS/MS spectra (**Figure 1a,b**). Using spectral libraries the tomatidine
187 molecular family was shown to contain both annotated nodes (level 2 or 3, according to the 2007
188 metabolomics standards initiative (Sumner et al., 2007) e.g, tomatidine, solasodine and
189 sarsasapogenin (**Figure 1b**), as well as unannotated nodes, which are also observed with
190 molecules occurring within Nightshade (Solanaceae) samples from the Global FoodOmics data
191 set (**Figure 1c**). Sarsasapogenin (**Figure 1c, node 1**) is found in food as well as stool data while
192 the +15.996 Da, the addition of the atom “O”, is only observed in stool data. However, numerous
193 other molecular families (such as **Figure 1c, node 10**) contain no annotation, but do have spectral
194 matches between plasma and foods — in this case features also observed in grape and fermented
195 grape samples. In other cases, a plasma metabolite is annotated and connected to unannotated
196 compounds found within the food reference samples (**Figure 1c, nodes 11-14**). These examples
197 highlight how molecular networking can be used to propagate potential metabolism. How potential
198 metabolism can be inferred with molecular networking is explained in (Quinn et al., 2017) and
199 (Aron et al., 2020).

200 A critical aspect of being able to leverage the food reference data, akin to expressed
201 sequence tags, is that the associated metadata can be retrieved and organized. We leverage the
202 Global FoodOmics ontology to identify different food categories in which MS/MS spectra are
203 observed. These food counts can be summarized for a dataset and then visualized as a flow chart
204 (**Figure 1d**). Due to the controlled research diets of the participants of the sleep and circadian
205 study in **Figure 1d**, we were able to report if a given food category was consumed during the
206 study. Of the 15 categories observed at level 5 of the food ontology, 8 represented direct matches,
207 3 represented fermented counterparts of consumed foods (such as yogurt and fermented grapes
208 when milk and grapes were consumed), and 4 categories were not documented to be consumed,
209 while coffee and tea were not provided to participants during this study. By and large, consistent
210 with the lack of consumption of caffeinated beverages, evidence of coffee or tea consumption

211 was only observed in two individuals. In one individual, caffeine was only detected in the first 48
 212 hrs, and in the other volunteer, caffeine was observed in a single time point in the later part of the
 213 study (second to last time point). Spectral matches to caffeine were not detected in any of the
 214 other participants. Thus, the empirically-recovered food ontology information from metabolomics
 215 data demonstrates that these matches are consistent with the food that was consumed in this
 216 study.
 217



218
 219 **Figure 1. The concept of interpretive metabolomics leveraging reference data sets.** a. A schematic
 220 overview of human data and reference data (e.g. data from food items) as molecular families from
 221 independent data sets that are used in b. b. A schematic representation when reference data is co-
 222 networked with human metabolomics data. Each node represents a unique MS/MS spectrum. c.
 223 Experimentally observed molecular families (sub-networks) generated from the co-analysis of stool (light
 224 blue) and plasma (magenta) data from a sleep restriction and circadian misalignment study with the Global
 225 FoodOmics reference dataset (green). Annotations are level 2/3 according to the 2007 metabolomics
 226 standards initiative (Sumner et al., 2007). Nodes 1-9: Tomatidine molecular family. Molecular family 10:
 227 a molecular family identified based on overlap of grape and fermented grape samples with plasma samples;
 228 multiple nodes contain spectral matches, however there is no library annotation and would otherwise remain
 229 completely uncharacterized. d. Summary of the spectra observed in plasma at each of the five food
 230 ontology levels. As this cohort received controlled diets, food categories observed in plasma samples were

231 matched with the known foods consumed. Solid circles represent MS/MS matches to foods consumed
232 during the study, while grey circles represent MS/MS matches to fermented versions of foods consumed,
233 indicating possible byproducts of digestion. Open white circle indicates consumption was not recorded in
234 this study.

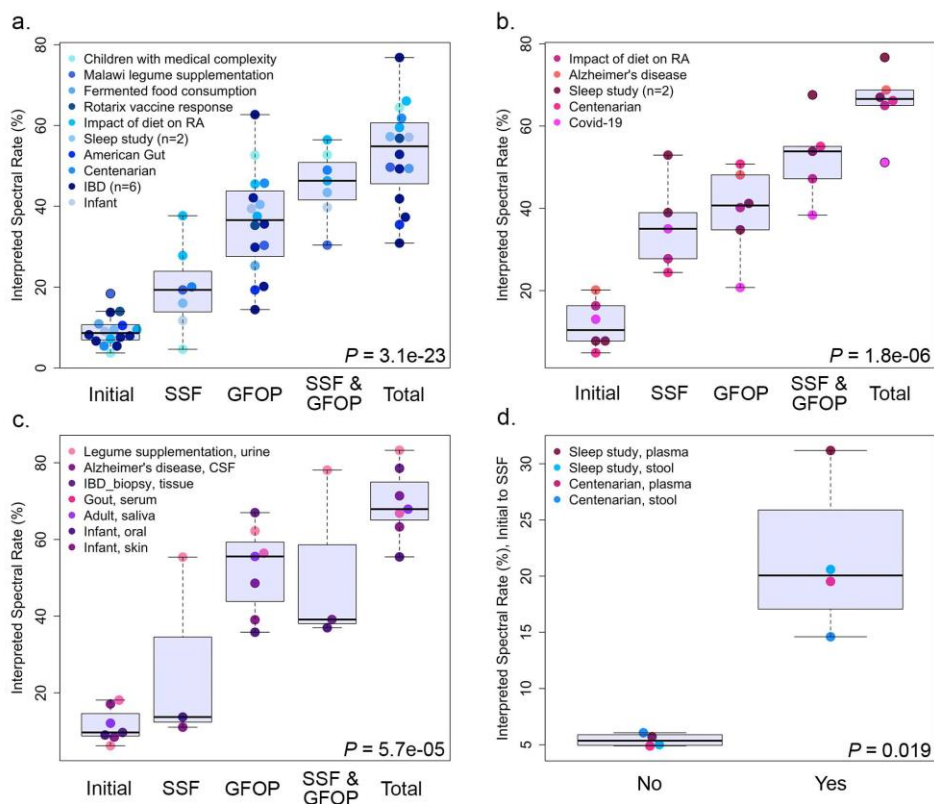
235
236 To illustrate the broad utility of the Global FoodOmics reference data in enhancing the
237 information gained from untargeted metabolomics, we co-analyzed the Global FoodOmics
238 dataset with 27 human datasets (**Table S2**; at 1% FDR spectral matching), with the inclusion of
239 additional study specific foods (SSF) where applicable (**Figure 1a**). These datasets contained
240 between 5 and 2123 samples, represented multiple different biofluids and tissues, and included
241 both adult and pediatric subjects, in conditions ranging from extremely long lived, such as a
242 centenarian-enriched population in the Cilento Blue Zone in Italy, to inflammatory bowel disease,
243 the healthy young adults undergoing experimental sleep restriction and circadian misalignment
244 highlighted in **Figure 1** (Sprecher et al. 2019), children with medical complexity, adults with
245 Alzheimer's disease, and Covid-19 infections in Brazil (**Table S2**).

246 Spectral matching to food reference data, observed as overlaps between datasets from
247 molecular networking, increased the interpretable fraction by 5.1 +/- 3.3 fold, even when
248 compared to the library of all 150,633 public reference spectra that are used by the GNPS analysis
249 infrastructure for annotation of public data which presently includes 29 spectral libraries, including
250 from the three MassBanks (Japan, EU and North America) (Horai et al., 2010), HMDB (Wishart
251 et al., 2018), ReSpect (Sawada et al., 2012), NIH natural product libraries (Huang et al., 2019),
252 PNNL lipid library (Kyle et al., 2017), Bruker/Sumner, FDA libraries, Gates Malaria library, EMBL
253 library, as well as many other GNPS contributed libraries
254 (<https://gnps.ucsd.edu/ProteoSAFe/libraries.jsp>) and the commercial NIST17 library (CID portion
255 only). Adding in additional information from molecular network connectivity, which can capture
256 metabolized versions of molecules, the fold change of interpretable data increased further to 6.8
257 +/- 3.5 fold (**Figure 2**). The Global FoodOmics reference samples significantly increased the
258 interpretation of various human metabolome samples above the initial annotation rate by 26.8+/-
259 3.3% for stool data ($P = 2.8e-16$, Games-Howell test), 27.5 +/- 5.2% for plasma data ($P = 0.0040$,
260 Games-Howell test) and 41 +/- 4.6% for other human data ($P = 0.00020$, Games-Howell test).
261 Further inclusion of connected nodes, representing potential metabolism via molecular
262 transformations, results in a total increase of 43.7 +/- 3.1% (fecal; $P = 6.9e-10$, Games-Howell
263 test), 51.2 +/- 6.9% (plasma; $P = 2.8e-06$, Games-Howell test), and 58.0 +/- 4.2% (human other;
264 $P = 1.4e-06$, Games-Howell test) percent of MS/MS spectra that can now be leveraged as
265 potentially a direct empirical readout of diet.

266 For 14 of the public datasets, food samples of the region or exact dietary items frequently
267 or exclusively eaten by that particular population were also collected (study specific foods; SSF).
268 SSF and Global FoodOmics reference samples were separately (SSF; GFOP) and jointly (SSF &
269 GFOP) evaluated for changes to the interpretable fraction of MS/MS spectra (**Figure 2**). For
270 example, adding SSF (n=38) alone increased the percent of interpreted spectra for the
271 centenarian stool data from an initial 5.4% annotation rate against spectral libraries to 20.0%
272 interpreted (**Figure 2a**) and 4.9% initial to 24.4% for plasma samples (**Figure 2b**), and adding
273 Global FoodOmics further expanded this to 49.0% (55.0% in plasma). For the sleep restriction
274 and circadian misalignment study highlighted in **Figure 1**, the interpreted fraction also increased
275 from an initial 7.2% to 27.8% (n=197 food samples; 45 of which are pooled meal samples), with

276 a further increase to 46.3% when using the Global FoodOmics reference data set (7.8% to 38.9%
 277 and with Global FoodOmics up to 54% for plasma). Overall, the inclusion of SSF significantly
 278 contributed to the increase in dietary spectral matches in plasma (**Figure 2b**; $P = 0.0028$, Games-
 279 Howell test). In addition, in some cohorts the interpreted spectral rate reaches almost 80% after
 280 expansion with molecular networking (**Figure 2c**).

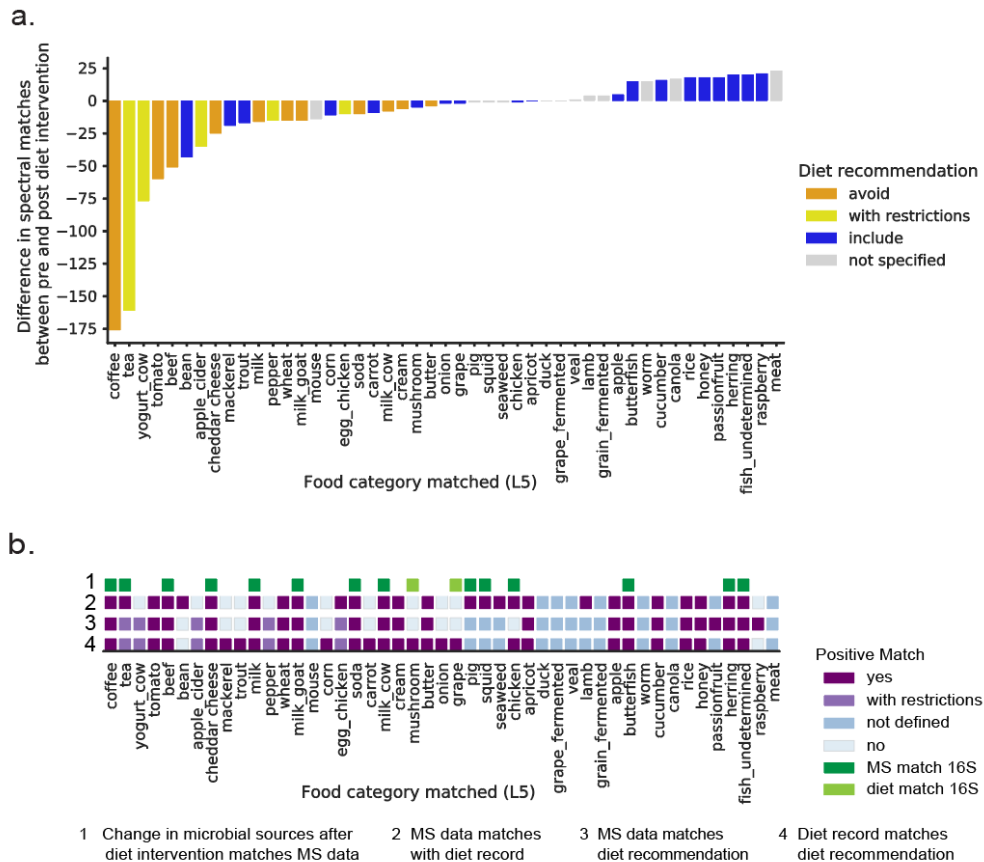
281 To further demonstrate that spectral matching using reference matching reflects dietary
 282 components, we performed a crossover study to test whether a mismatched SSF inventory would
 283 yield similar results to the increases observed across studies with SSF (e.g. centenarian foods
 284 for the sleep and circadian study cohort). Crossover revealed that the reciprocal tests
 285 interpretation rates were only a few percent (5–6%) in comparison to when the correct SSF were
 286 used (15–30%) (**Figure 2d**).



287
 288 **Figure 2. Increases of MS/MS spectral match rates when using interpretive metabolomics at the**
 289 **data set level.** Spectral match rates of molecular features due to library match, food reference data, and
 290 molecular networking in **a.** stool data. Significant differences are determined by Welch's F-Test. Library
 291 spectral matches (initial), spectral matches to study specific foods (SSF), spectral matches to Global
 292 FoodOmics project (GFOP) data, both (SSF & GFOP), expansion with molecular networking (Total). **b.**
 293 plasma data, and **c.** other human biospecimens. **d.** A crossover experiment between the centenarian data
 294 from Italy and the sleep and circadian study from the US, for both fecal and plasma samples. Study specific
 295 foods consumed by those individuals (yes) vs a different set of study specific foods (no), (Welch's *t*-test).
 296

297 As the Global FoodOmics reference database expands with regionally-specific foods
 298 through a continued community effort, the interpreted fraction will likely increase. For example,
 299 when legume food data (15 files; SSF) similar to legumes supplemented in an infant malnutrition

300 study were included in addition to the Global FoodOmics data, the number of spectral counts for
 301 legumes went from 105 to 2430 unique MS/MS spectra that matched, while other food categories
 302 such as dairy and meats remained constant (level 3 food ontology; Legume supplementation,
 303 urine). Regional specificity was also directly evident for plasma samples collected in Brazil for a
 304 Covid-19 study, which displayed more spectral matches to a locally collected set of 60 Brazilian
 305 food samples with ~35% increase than to the entire Global FoodOmics reference dataset, that is
 306 dominated by US food, which only gave an ~20% increase in spectral matches (**Figure 2b**). Thus,
 307 although there is some overlap among the data from different foods, and even overlap among
 308 human-derived metabolites and the food data (e.g. many primary metabolites or those common
 309 in vertebrates), a large proportion are sufficiently unique to reveal, at least in part, the dietary
 310 composition in the study.



311 **Figure 3. Using interpretive metabolomics to assess dietary recommendations at the study level. a.**
 312 Food ontologies of a rheumatoid arthritis cohort before and after a specific dietary recommendation of a
 313 low inflammatory diet. Plasma data are used. Food categories indicated as 'with restrictions' encompass
 314 foods where different types are encouraged and others discouraged (green vs. black tea) or foods that were
 315 supposed to be minimized (such as limiting egg consumption to 2 eggs per week). Food categories
 316 indicated as 'not specified' could not be matched to the suggested diet. **b.** Comparison of interpretive
 317 metabolomics results in the recommended diet and self-reported diet intake. Diet diaries were tabulated as
 318 consumption or no consumption of >200 food categories over the 28 days of the study and matched to the
 319 MS food categories, as possible. Matches to 16S rRNA gene sequence data are based on Bayesian source
 320 tracking proportions from the bacterial community, with food types as sources and rheumatoid arthritis fecal
 321 samples as sinks. The increase or decrease in the proportion of food source contribution pre and post
 322 dietary intervention (y-axis) is colored according to dietary recommendations.
 323

324

325 To assess if interpretive metabolomics could be used to empirically establish adherence
326 to dietary recommendations using MS/MS data, we analyzed a data set from rheumatoid arthritis
327 patients (RA) asked to follow an anti-inflammatory diet (ITIS diet) (Bustamante et al., 2020). We
328 compared the per sample extracted food counts with the recommended diet alteration as well as
329 self-reported diet diary entries. The recommended diet included some foods to be avoided (such
330 as coffee, refined sugars and milk), some foods to be restricted (minimize red meat and egg
331 consumption) and some foods to be frequently consumed (such as fruits/vegetables, and plain
332 unsweetened yogurt). In total, 47 foods and beverages were observed in this project with
333 interpretive metabolomics (**Figure 3a**). By and large, most adhered to the recommended diet, as
334 food counts of recommended foods increased, and those of foods to avoid decreased. Although
335 there are instances when the mass spectrometry based observations did not match the
336 recommended diet regime, the self-reported dietary records matched the empirically determined
337 foods better than the recommended dietary information (**Figure 3b**). We further validated these
338 matches using source tracking with 16S rRNA gene amplicon data collected on ~1500 samples
339 of the Global FoodOmics foods, to predict food source contribution to the RA study stool samples.
340 We observed a highly significant correlation in the proportion change of food sources predicted in
341 the stool samples and metabolites in the plasma before and after dietary intervention (Pearson
342 $r = 0.57$, p -value = 0.003; **Figure S2**). The empirically recovered food ontology information from
343 interpretive metabolomics, in conjunction with validation with DNA sequence data, demonstrates
344 the ability to recapitulate dietary readouts from human biospecimens and assess diet adherence.

345 Interpretive metabolomics comes with several caveats to consider. We are not yet able to
346 capture a complete picture of the human diet: for example, in the RA study, the participant diet
347 diaries contained foods not yet captured in the FoodOmics database, potentially leading to an
348 underestimation of food types observed. Community expansion of the Global FoodOmics
349 database with specific foods and food ingredients will ultimately eliminate this issue.

350 Another consideration is similar to what is observed with expressed sequence
351 tags/RNASeq, where it is common to observe that there are multiple sample types, tissue
352 locations or conditions that result in misinterpretation because the same sequence occurs in
353 multiple locations. By analogy, a molecule could be produced by humans but also be part of
354 different diet sources (i.e. cholesterol produced by the human body versus consumed). However,
355 such matches still enable one to formulate a hypothesis that the observed MS/MS features from
356 the human data might originate from the reference data as source, in this case food, especially
357 when there are hundreds or thousands of signatures that point to specific foods or food groups
358 that overlap.

359 As we saw in many of the above datasets, it is not atypical to observe small numbers of
360 spectral matches to insects, rodents, fungi and worms within diet read-outs. Although data on
361 fungi, tarantula, crickets, and black ants, meant for human consumption, are included, most of
362 these samples that match human data sets are from a Global FoodOmics sampling effort at the
363 San Diego Zoo. While there is likely some overlap with molecules from these less common foods
364 to those that humans more commonly consume (e.g. certain acylcarnitines might be found in beef
365 and mice), the FDA food contamination guidelines allow for insect, fungal, worm, rodent parts and
366 fecal matter to be present in food in quantities that surprise many non-specialists (Center for Food
367 Safety and Nutrition, 2019) For example, peanut butter is allowed to have 30 or more insect

368 fragments and one rodent hair per 100 grams, and apple butter is allowed to have “5 or more
369 whole or equivalent insects (not counting mites, aphids, thrips, or scale insects) per 100 grams of
370 apple butter.” As long as these dietary “additives” are added to the reference data set, they too
371 will be observed. Thus, interpretive metabolomics can provide empirical support for dietary
372 compliance in nutritional content, including in clinical studies, and capture information that would
373 otherwise remain hidden.

374 Conclusion

375 Here we show that well-curated reference datasets can be leveraged to provide a deeper
376 understanding of untargeted metabolomics. Adding food-based spectral matches improves our
377 ability to interpret molecular features 2 to 14-fold, and further improves to 3 to 17-fold by
378 incorporating connections from molecular networking, providing a deeper insight of the
379 metabolomic “dark matter.” Our results indicate that a direct empirical readout of diet adherence
380 is within our reach using interpretive metabolomics, by combining structural, source, and chemical
381 similarity measures.

382 Although we demonstrated the power of interpretive metabolomics with food data as
383 reference, any individual reference data set or combination of multiple data sets could be used in
384 this fashion. We envision the broad application of such an approach. Generating databases for
385 environmental allergens, medications, illegal substances, food ingredients and personal care
386 products can inform within those research areas on potential exposures and food adulteration.
387 Further, such investigations may also have far reaching impacts to understand commonalities
388 that underlie different diseases. Over time, as metabolomics data repositories begin to control
389 metadata vocabularies, most public data could be leveraged and reused as a reference data set
390 on its own. This will significantly improve the interpretability of all metabolomics data, be it from
391 environmental, animal, or human sources.

392 Acknowledgments

393 Funding sources: We thanks the CCF foundation #675191, U19 AG063744 01, R01AG061066,
394 1 DP1 AT010885, P30 DK120515 Office of Naval Research MURI grant N00014-15-1-2809 and
395 NIH/NCATS Colorado CTSA Grant UL1TR002535. This work was also supported in part by the
396 Chancellor’s Initiative in the Microbiome and Microbial Sciences and by Illumina, Inc. through
397 reagent donation in partnership with the Center for Microbiome Innovation at UC San Diego. We
398 would like to thank Elaine Wolfe and Karenina Sanders for sample processing, and Jeff DeReus
399 for data handling, processing and maintaining the computational infrastructure. JPS was
400 supported by SD IRACDA (5K12GM068524-17), and in part by USDA-NIFA (2019-67013-29137)
401 and the Einstein Institute GOLD project (R01MD011389). RC and MG were supported by Krupp
402 Endowed Fund. RHM was supported through a UCSD training grant from the NIH/NIDDK
403 Gastroenterology Training Program (T32 DK007202). The Brazilian National Council for Scientific
404 and Technological Development (CNPq)-Brazil [245954/2012] to MFO. KS was supported by a
405 PROMOS fund (DAAD). WB is a postdoctoral researcher of the Research Foundation – Flanders
406 (FWO). We thank Ricardo da Silva for his feedback and early bioinformatics analysis for the
407 Global FoodOmics project. We further acknowledge all the individuals that contributed samples

408 as well as companies and organizations that have donated samples: Townshend's Tea Company,
409 BDK Kombucha, Oregonian Tonic, Squirrel & Crow, Venissimo cheese, Fermenter's Club San
410 Diego, Good Neighbor Gardens, Sprouts Farmers Market, Ralphs, Whole Foods and SD Zoo and
411 Safari Park. Specifically thank you to Austin Durant for coordinating sampling at Fermentation
412 Festivals and the wonderful staff at SD Zoo Global for coordinating and helping with sample
413 collection: Michele Gaffney, Edith Galindo, Katie Kerr, Andrea Fidgett, Jennifer Stuart, Debbie
414 Tanciatco, and Lisa Pospychala.

415 Author Contributions

416 PCD, RK, RJD, and JMG conceptualized the idea.
417 MWP, FDO, KCW, CMA, EB, KS, PCD, RJD, RK, NCS, ADS, GA, DM, NPL, and JMG collected
418 FoodOmics samples and performed metadata curation.
419 MWP, FDO, FV, CMA, EB, NCS, and JMG performed FoodOmics sample processing and MS
420 data acquisition.
421 AJJ, PBF, ED, QZ, DN, DM, JPS, and JMG curated Global FoodOmics metadata to match
422 FNDDS.
423 JBW, BSB, BJB, RC, MGDB, MD, EOE, DG, LH, JK, MM, CM, RK, KES, DVR, CW, KPW, MFO,
424 RHM, DW, RT, JGA, PD, MG, DG, AKJ, BJB, RMS, KCW, ADS, FV, NPL, and JMG provided
425 samples, comparative dataset, and/or detailed metadata.
426 LMMM, TMC performed Covid-19 patient and/or food sample preparation and analysis.
427 PLJ was the physician responsible for the Covid-19 patients (provided samples).
428 RDRO was the physician responsible for collecting the plasma from Covid-19 patients (provided
429 samples).
430 FPV was responsible for tabulation of Covid-19 patient data and analysis of health components,
431 and body dimensions.
432 TS, MB, LDG, GH performed 16S sequencing and prep.
433 CM, DM, JPS performed source tracking and/or 16S data analysis.
434 MW supported GNPS computational infrastructure used in the study.
435 CLW, WB, AKJ, ES, AT, NPL and JMG analyzed MS data.
436 CLW, WB, AKJ, CM, and JMG generated figures.
437 PCD, RK, RJD, ADS, and JMG supervised the work.
438 PCD, RK, CLW, and JMG wrote the paper.
439 All authors have contributed feedback and edits to the manuscript.

440 Declaration of Interests

441 BSB has a research grant from Prometheus Biosciences and has received consulting fees from
442 Pfizer. PCD is on the scientific advisory board of Sirenas, Cybele Microbiome, Galileo and founder
443 and scientific advisor of Ometa Labs LLC (with approval by UC San Diego). JHK is a consultant
444 for Medela, Astarte Medical, Nutricia, and Fujifilm; he owns shares in Astarte Medical and
445 Nicolette. MG has research grants from Pfizer and Novartis. PSD has received research support
446 and/or consulting from Takeda, Pfizer, Abbvie, Janssen, Prometheus, Buhlmann, Polymedco.
447 AJJ has received consulting fees from Abbott Nutrition and Corebiome. DG is a consultant for

448 Biogen, Fujirebio, vTv Therapeutics, Esai and Amprion and serves on a DSMB for Cognition
449 Therapeutics. KPW reports during the conduct of the study receiving research support from
450 SomaLogic, Inc., consulting fees from or served as a paid member of scientific advisory boards
451 for the Sleep Disorders Research Advisory Board - National Heart, Lung and Blood Institute,
452 CurAegis Technologies, Circadian Therapeutics, LTD. and Circadian Biotherapies Ltd. ADS and
453 RK are directors at the Center for Microbiome Innovation at UC San Diego, which receives
454 industry research funding for multiple microbiome initiatives, but no industry funding was provided
455 for this project. MW is a co-founder of Ometa Labs LLC.

456 STAR Methods

457 Resource Availability

458 Lead Contact

459 Further information and requests for resources should be directed to and will be fulfilled by the
460 Lead Contact, Pieter Dorrestein (pdorrestein@health.ucsd.edu).

461 Materials Availability

462 This study did not generate new unique reagents.

463 Data and Code Availability

464 The code generated during this study is available on GitHub at
465 <https://github.com/DorresteinLaboratory/GlobalFoodomics>.

466 Raw and processed 16S rRNA amplicon sequencing data is available at Qiita study
467 #11442 and raw sequence data has been deposited at EBI accession ERP122648.

468 GNPS task ID of analysis used for tag generation: f1a1f3a61aca416a9b3687d72488da7f

469 The following files are available in addition to the Global FoodOmics mzXML files on
470 massive.ucsd.edu under MSV000084900: metadata as a .txt; an image repository with between
471 1 and 6 images per food; table of FDR-based parameters; raw food count data for RA study; full
472 size PDF of sleep restriction and circadian misalignment study - GFOP3500 molecular network
473 (excerpts found in Figure 1).

474 Metadata dictionary:

475 [https://docs.google.com/spreadsheets/d/1Ebn-](https://docs.google.com/spreadsheets/d/1Ebn-TgMWEkd_7KOW9TCRvHGPsE7dGjVCr7dq28pwbmM/edit#gid=727944641)
476 [TgMWEkd_7KOW9TCRvHGPsE7dGjVCr7dq28pwbmM/edit#gid=727944641](https://docs.google.com/spreadsheets/d/1Ebn-TgMWEkd_7KOW9TCRvHGPsE7dGjVCr7dq28pwbmM/edit#gid=727944641)

477

478 The GNPS analyses used in this study can be accessed on-line at the following links:

- 479 • Sleep study (MSV000083759;
480 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=e0bf255bcb2e492bb0be3be1a691b5fb>;
481 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=6fe434761daf4f9da540cf1fd90b3985>;
482 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=9a90bd12f51e453e968656e6458e0da4>)

- 483 • Centenarian (MSV000084591;
484 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=8895b6e3445546c4a5bc3a726a920227>;
- 485 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=981c9a7d39f742bda296d52f856981e5>)
- 486 • Impact of diet on RA (MSV000084556;
487 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=0794151fce2c4c18a7a0aa3a09140169>)
- 488 • LP Infant (MSV000083462; MSV000083463;
489 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=a7b222466ef844e69cdbc9835d2f6c39>;
- 490 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=c756a9dfb5c34a2a8655f88114edf0a8>;
- 491 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=4a322e640bb644068030949267fb4ea9>)
- 492 • Children with Medical Complexity (MSV000084610;
493 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=df24423835a341969342c2086b46275a>)
- 494 • American Gut (MSV000081981;
495 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=4884483bcffe4f269819858c3fd4faef>)
- 496 • Fermented food consumption (MSV000081171;
497 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=5cca39e0ebab4066a56e41ded48b4466>)
- 498 • Malawi legume supplement (MSV000081486;
499 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=93ba727aa9234727a73ae7860b2af3ca>)
- 500 • Rotarix vaccine response (MSV000084218;
501 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=08e9b9e048f04ac4b416e574a073e8e6>)
- 502 • IBD_1 (MSV000082431;
503 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=ec08eed8f186430d893c63111409baf4>)
- 504 • IBD_individual (MSV000079115;
505 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=fad746939afd4184975a296436aebfb7>)
- 506 • IBD_seed (MSV000082221;
507 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=907f2e0b7878417dbdb4c83f0df0e83a>)
- 508 • IBD_biobank (MSV000079777;
509 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=a79fbd4c96124209adfd0ef84cb56dec>)
- 510 • IBD_2 (MSV000084775;
511 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=07f855658c5342458045032ea70fc526>)
- 512 • IBD_200 (MSV000084908;
513 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=55bef02250d744eb97c6040c379cbfb4>)
- 514 • Alzheimer's disease (MSV000085256;
515 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=aac78e9d23b84194ab2f768cb685c636>)
- 516 • Covid-19 (MSV000085505; MSV000085537;
517 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=9cbcb6b46fe24826bc56c9e893d0bd2b>)
- 518 • IBD_biopsy (MSV000082220;
519 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=a83a279dad154f9ca7b549d40ce117ba>)
- 520 • Gout (MSV000084908;
521 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=55bef02250d744eb97c6040c379cbfb4>)
- 522 • Adult Saliva (MSV000083049;
523 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=6dd6e5b1cf454d67b8a2b3c151c18f4a>)
- 524 • Legume supplementation (MSV000084663;
525 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=93ba727aa9234727a73ae7860b2af3ca>)
- 526

527 Networking parameters were set based on the MOLECULAR-LIBRARYSEARCH-FDR workflow
528 on GNPS with the following task IDs:

- 529 • GFOP3500: a7bf6cc3f91d466bab923f2268d6f4fc
- 530 • Sleep deprivation: b55ab4004ed342d7b4ed1c488e935998
- 531 • Sleep study: 78bbfed8574748d1a77dc7c2f1a44d39
- 532 • Sleep study_SSF_test: b55ab4004ed342d7b4ed1c488e935998
- 533 • Centenarian: 265a9553c69e47499cca3de056b43178
- 534 • Centenarian_SSF_test: 265a9553c69e47499cca3de056b43178
- 535 • American Gut: aee5dde3b2f84079a264e68ec981487e
- 536 • Fermented food consumption: a44d1b2e1b9d4612974d0b85021675a7
- 537 • Malawi legume supplement: de7b55f8adaa4ad9b2a8430e30435bf3
- 538 • Children with Medical Complexity: f27243af071b43ab90d846bda959fc1c
- 539 • Rotarix vaccine response: a2e02e3f97a54ca08e3866cc60f8d42b
- 540 • Impact of diet on RA: 62b8754e761549f3b94ffae83d7ab95a
- 541 • LP infant: 532aba2ad3644fadba0e6e7ea063c7ee
- 542 • IBD_1: bb10b1ce90a24f3a9cef1e85e88c3882
- 543 • IBD_biopsy: c4cfda90933b4842a7154f5f2def139d
- 544 • IBD_individual: 3ce8cc636ae944848b4ada322aaf12fe
- 545 • IBD_seed: ebbb715fc605457ba5f7e910b79d6177
- 546 • IBD_biobank: 9465c34cf5444e12b89318b1fb363714
- 547 • IBD_2: 983fa9271136404fb5743b44a6a109f0
- 548 • IBD_200: e5acf5726722486caa897b2b07d402e8
- 549 • Impact of diet on RA: 62b8754e761549f3b94ffae83d7ab95a
- 550 • Alzheimer's disease: 658103164325425981c097cecba840b0
- 551 • Gout: a478f419ae824378aa02e5e1b310cad2
- 552 • Adulta saliva: 32980f95dbd5437aaa9e15d05c7246bb
- 553 • LP infant: 8fbfbc1bf38c418fb223306cd42af897
- 554 • LP infant: 3e414e13a4394bb78c07f7ca7f4d1be3
- 555 • Legume supplementation: 2ca007303b9c4bb3820f392b996eba27
- 556 • Alzheimer's disease: 658103164325425981c097cecba840b0
- 557 • Covid-19 Brazil: d16eb32276c84bdb9c35c5872e97a986

558 Methods

559 IRB information for the human datasets used in this study, made public on MassIVE
560 Sleep study (MSV000083759; IRB 15-0282), centenarian (MSV000084591; IRB 180478), Impact
561 of diet on RA (MSV000084556; IRB 161474), LP Infant (MSV000083462; MSV000083463; IRB
562 151713 UCSD), Children with Medical Complexity (MSV000084610; IRB 161948 UCSD),
563 American Gut (MSV000081981; IRB 141853 UCSD), Fermented food consumption
564 (MSV000081171; IRB 141853 UCSD / published), Malawi legume supplement (MSV000081486;
565 IRB ID #201503171; Washington University Human Studies Committee), Rotarix vaccine
566 response (MSV000084218; IRB is PR-10060 from University of Virginia), IBD_1
567 (MSV000082431; IRB # 150675), IBD_individual (MSV000079115; IRB # 150675), IBD_seed

568 (MSV000082221; UCSD HRRP 131487), IBD_biobank (MSV000079777; UCSD HRRP 131487);
569 IBD_2 (MSV000084775; IRB # 150675), IBD_200 (MSV000084908; IRB # 150675), Alzheimer's
570 disease (MSV000085256; UCSD IRB # 170957), Covid-19 (MSV000085505; MSV000085537;
571 IRB approval number is 30248420.9.0000.5440 (University of São Paulo, Brazil), IBD_biopsy
572 (MSV000082220; IRB number is 120025), Gout (MSV000084908; IRB Project #160768X), Adult
573 Saliva (MSV000083049; IRB 150275 UCSD), Legume supplementation (MSV000084663; IRB ID
574 #201905103).

575 Global FoodOmics

576 Sample Collection

577 Sampling methodology was developed in order to facilitate sample collection in any environment,
578 from the home, a restaurant, a festival, or in the lab. Initial samples were collected in a consistent
579 manner, between April 2017 and March 2018. Additional sets of samples were added through
580 Fall 2019. Each sample was assigned a unique number identifier upon sampling, which was used
581 to trace the origin of the sample, and to organize descriptive information about the sample. In
582 addition, when possible samples were photographed by the participant to create a photographic
583 archive of all samples (uploaded to MassIVE MSV000084900; >4000 images representing 67%
584 of the samples (2399/3579)). Primarily for the initial data set these images were used as the first
585 point of reference for the collection of ancillary information about the different samples (termed
586 metadata, described in more detail below). The image archive was critical, because as the project
587 evolved and the breadth of sample types increased, new categories were added to the metadata,
588 which were then filled in weeks or even months after sample collection.

589 Samples were frozen at -80°C within 24 h of sample collection, unless otherwise noted in
590 the metadata. Two samples were collected for each food or beverage included in the study. One
591 sample was collected as an archive and directly frozen, and a second sample was collected for
592 extraction. Food samples were collected in a tube prefilled with 1 ml 95% ethanol (Ethyl alcohol
593 (Sigma- Aldrich) and Invitrogen UltraPure™ Distilled Water), as high ethanol concentrations are
594 efficacious at preserving the sample for both DNA and metabolite analyses (Song et al., 2016).
595 Samples were collected into 2 ml round bottom microcentrifuge tubes (Qiagen) and weighed prior
596 to freezing. The pre-sample and post-sample weights as well as the weight differences were
597 recorded in the metadata. It was not possible to collect all samples at a given concentration of
598 extraction solvent (ethanol), because sampling was performed in many different environments
599 and is meant to be consistent with future crowd-based community science participation.
600 Therefore, the data can be compared qualitatively and not quantitatively, however for certain
601 subsets 50 mg of material was collected consistently.

602 Additional sets of food samples were added to the core set using the same methods as
603 outlined above when possible. Samples from Venezuela were collected whole in absolute ethanol
604 $\geq 99.8\%$ (Sigma Aldrich) and the extract was processed directly.

605 The experimental protocol for the sleep restriction and circadian misalignment study has
606 been described previously (Sprecher et al., 2019). Meals and food samples were prepared by the
607 Clinical and Translational Research Center Nutrition Core of the Colorado Clinical and
608 Translational Sciences Institute. Food was transported to the research site and refrigerated for
609 the duration of the in-patient study. Individual meals were sampled and stored frozen in ziploc

610 bags. They were stored at -80°C prior to subsampling and LC-MS/MS analysis. Images are
611 contained in a separate Sleep Study folder (MSV000084900).

612 For several of the human studies we collected data on associated foods, which were
613 processed according to the same methods as the Global FoodOmics samples. The number of
614 SSF samples per cohort are outlined here: experimental sleep deprivation (197 samples; 45 are
615 pooled); centenarian (38 individual samples); malawi legume supplement (14; 2 sample types,
616 several extraction types); children with medical complexity (24 formula samples; 11 exact
617 overlap); RA diet samples (20 individual sample; 2 samples types (stool, plasma), 2 time points));
618 mother's milk (58 milk samples); legume supplements (15 individual legume samples; 6 different
619 types).

620 Community-based science collection

621 The first sample collected was a carrot from a home garden. The participant was interested in
622 how the soil conditions from prior tenants would impact the chemistry of the carrot, since the
623 gardening practices of the prior tenant were unknown (organic or not, pesticide usage, etc.). In
624 addition, home grown foods often taste different than store bought, likely reflected in the food
625 metabolome.

626 During the course of sampling, samples were received from over 50 different individuals
627 in California as well as from different states as well as countries (such as Venezuela, Italy and
628 most recently Brazil). Contributions from individuals ranged from produce from home gardens,
629 home fermented products (yogurt, kombucha, sauerkraut), meat and dairy from private farms, to
630 items individuals had purchased that were of interest to them.

631 We were also directly invited to sample at local stores and organizations, including
632 Venissimo cheese, Good Neighbor Gardens, and the San Diego Zoo and San Diego Zoo Safari
633 Park, as well as local supermarkets such as Sprouts Farmers Market, Whole Foods Market, and
634 Ralphs. We were invited by San Diego Fermenter's Club founder Austin Durant to the San Diego
635 Fermenter's Club meeting and sampled from multiple vendors at both the Oregon Fermentation
636 Festival in 2017 as well as the San Diego Fermentation Festival in 2018. We also received citrus
637 samples from a farm at the US-Mexico border, with visibly dark skin due to air pollution, a
638 particular concern of the farmer. Other sampling occurred in conjunction with study design, as
639 was the case for the Rheumatoid arthritis cohort and the Covid-19 study. In total we engaged with
640 a broad range of individuals, organizations, businesses and scientists, to generate this dataset of
641 3579 samples (for future use this is already expanded beyond this number due to the collection
642 of sets of SSF). A predominance of foods included in this initial dataset were sampled and/or
643 purchased in California, leaving room for much further expansion and the inclusion of a crowd-
644 sourced community science initiative to expand the array of samples.

645 The sample set contains a broad set of simple foods including fruits, vegetables, grains,
646 as well as raw meat and fish, which build the foundation of many food products. In addition, we
647 have 1133 fermented samples. This subcategorization of foods is made possible by the metadata
648 collected on these samples, described in the Metadata Curation section. The breadth of samples
649 included in the dataset necessitated careful collation and a range of information about the
650 samples, resulting in 157 different metadata categories to describe various aspects of these food
651 and beverage samples.

652 Samples originate from over 50 different identified countries of origin (Argentina, Australia,
653 Austria, Belgium, Bolivia, Brazil, Canada, Chile, China, Columbia, Croatia, Ecuador, England,
654 Ethiopia, France, Germany, Greece, Guatemala, Haiti, Holland/Netherlands, India, Indonesia,
655 Ireland, Israel, Italy/Sardinia, Japan, Kenya, Korea, Madagascar, Malawi, Mexico, New Zealand,
656 Nilgiri, Peru, Philippines, Poland, Serbia, Portugal, Russia, Scotland, South Africa, Spain,
657 Switzerland, Taiwan, Thailand, Trinidad & Tobago, Turkey, UK, USA/Puerto Rico, Vietnam,
658 Venezuela; EU, South America not included separately).

659 Metadata Curation

660 General organization

661 Detailed information about each sample was captured in the form of metadata. The metadata are
662 in the form of an array, where each row represents one sample and each column captures unique
663 information about the sample (See Supplementary Information for Metadata File, as well as
664 updates on Massive MSV000084900). This matrix allows for the categorization of foods by
665 various different attributes and links these attributes to the sample numbers, the data files
666 (.mzXML filename), as well as the 16S sequence information on Qiita (sample_name). The initial
667 metadata categories captured included sample description, sample number, location sample was
668 collected, the weight of the sample (pre-sample, post-sample, sample weight), the day it was
669 collected, and whether an image had been taken and renamed to match the sample number and
670 archived in the image repository. The initial 9 categories captured minimal information and
671 allowed tracking of information about the sample.

672 During the process of sample collection, the diversity of the samples being collected
673 necessitated the addition of columns to capture more information about the samples and to be
674 able to categorize them and compare different attributes. These columns grew to capture highly
675 detailed information about each sample, for example whether the sample was organic, if it was
676 raw or cooked, if it was washed before sampling, or for cheese samples whether it is the rind or
677 the curd, etc. As columns were added, the initial columns and the image repository were used to
678 trace back information.

679 *Classification scheme*

680 Various classifiers are used to describe foods, however we were unable to find an established
681 scheme able to capture the diversity of samples, as well as distill the metadata down into a
682 manageable number of categories to distinguish differences between the metabolomes of
683 different food classes. We therefore categorized the foods by sample_type, which captured
684 whether the sample was a food, beverage, or other item (for example supplements) and then
685 expanded and shaped a unique categorization which takes into account the species and botanical
686 definitions of foods. The sample_type categories range from sample_type_land_aquatic, to
687 differentiate items sourced from different physical environments, sample_type_common, which
688 allows for representation of a particular food group which was not otherwise captured in the
689 metadata, such as zoo food or candy. The sample_type groups also include a hierarchy from
690 group1 to group6 (Levels 1 through 5 are referenced in this manuscript), specific to foods and
691 groupB1 through groupB3 which contain beverage specific information (alcoholic [binary],
692 carbonated [binary], type of beverage [such as red wine, kefir, soda, etc.]).

693 *Complex samples*

694 The above classification scheme gave sufficiently detailed information about simple foods (ones
695 that have only one ingredient and could thus be filled out to the last group level, such as red cherry
696 tomato). Complex foods contain not only multiple ingredients, but include highly processed foods
697 purchased with ingredient lists as well as home cooked or restaurant meals. These foods have a
698 higher variability of information known about them. The top 6 ingredients are captured in individual
699 metadata categories, with a seventh ingredient field which contains the remainder of the
700 ingredients (if known). However, the order of ingredients does not always clearly reflect the type
701 of food and some constituents that may be of interest, such as tree nuts which may only be found
702 in trace quantities. The `sample_type_common` category captured some of the information about
703 the type of sample (candy), however to have a tangible classification of different ingredient types,
704 we generated a specific complex food ontology based on the known presence of common
705 categories (corn, dairy*, egg*, fruit, fungi, fish*, shellfish*, meat, peanut*, seaweed, soy*, tree nut*,
706 vegetable/herb, wheat* (*designates known food allergen)). These categories reflect the main
707 food groups and some of the most common allergens (US FDA Food Allergen Labeling And
708 Consumer Protection Act of 2004) (Sicherer and Sampson, 2006), items which are of interest
709 when correlating food metabolome data with other datasets, such as human fecal material (where
710 the foods eaten are known or unknown).

711 *Fermented foods*

712 Preservation and processing method are included in the metadata. However, due to the potential
713 importance of fermentation in the alteration of the food metabolome, and the potential health
714 benefits that have been ascribed to fermented foods, several categories were included to highlight
715 this feature: fermented or not, whether it contains live active cultures, whether it contains
716 chocolate (which then was cross checked with the fermented category, as chocolate is a
717 fermented food). The list of fermented foods crosses many of our sample types as it includes
718 fermented dairy (yogurt, cheese), fermented meat/fish (salami, fish sauce), fermented vegetables
719 (kimchi, sauerkraut), fermented fruit (chocolate, coffee), and fermented grains/legumes (bread,
720 tempeh).

721 *Food specific categories*

722 Certain individual food categories also necessitated creation of specific categorization. For
723 example, cheeses have the specific categories `cheese_part` (curd vs. rind), `cheese_type`
724 (washed, blue, etc), and `cheese_texture` (soft, semi-soft, semi-hard, hard). Particularly for raw
725 plant products, such as fruits, vegetables, grains which form the basis for many food ingredients,
726 we captured botanical information: `botanical_anatomy` (fruit, leaf, tuber, seed, etc.),
727 `botanical_genus`, and `botanical_genus_species` (when known). Tea samples have tea quality and
728 tea type as distinct categories.

729 Metadata for Cross-study Comparison

730 To facilitate cross study comparison, we included the Earth Microbiome Project ontology: `emp_o_1`
731 (level 1: Free-living, Host-associated, Control, or Unknown), `emp_o_2` (level 2: Saline, Non-saline,
732 Animal, Plant, or Fungus), and `emp_o_3` (level 3: most specific habitat name)
733 [http://www.earthmicrobiome.org/protocols-and-standards/emp_o/]. Wherever possible we linked

734 foods to food identifiers or created identifiers and categories that built upon the existing framework
735 as defined by the U.S. Department of Agriculture's Food and Nutrient Database for Dietary
736 Studies 2011-2012 (FNDDS) food grouping scheme (Martin et al., 2012).

737 Sample Preparation

738 A sterile stainless steel bead was added to each sample collected in ethanol and the samples
739 were thawed on ice for 30 min. Samples were homogenized at 25–30 Hz for 5 min using a tissue
740 homogenizer (QIAGEN TissueLyzer II, Hilden, Germany). Samples were swabbed with sterile
741 dual tip swabs (BD swubes) and frozen immediately at -80°C until DNA extraction.

742 DNA Extraction and 16S rRNA gene amplicon sequencing

743 DNA extraction and 16S rRNA gene amplicon sequencing were performed using Earth
744 Microbiome Project (EMP) standard protocols (<http://www.earthmicrobiome.org/protocols-and-standards/16s>) (Thompson et al., 2017). DNA was extracted with the Qiagen MagAttract
746 PowerSoil DNA kit as previously described (Marotz et al 2017). Amplicon PCR was performed on
747 the V4 region of the 16S rRNA gene using the primer pair 515f-806r with Golay error-correcting
748 barcodes on the reverse primer. Amplicons were barcoded and pooled in equal concentrations
749 for sequencing. The amplicon pool was purified with the MO BIO UltraClean PCR cleanup kit and
750 sequenced on the Illumina MiSeq sequencing platform. Raw sequence data were uploaded to
751 Qiita for pre-analysis processing (Qiita study ID: 11442) (Gonzalez et al., 2018). In Qiita, raw
752 sequence data were demultiplexed and minimally quality-filtered using the QIIME 1.9.1 script
753 `split_libraries_fastq.py`, with a Phred quality threshold of 3, allowing for reverse complemented
754 barcodes and mapping barcodes, and default parameters. Demultiplexed, quality-filtered
755 sequence data were then trimmed to a read length of 150-bp, denoised with Deblur v1.1.0 (Amir
756 et al., 2017) using default parameters, and subject to fragment insertion with SATÉ-enabled
757 Phylogenetic Placement (Janssen et al., 2018) into the GreenGenes 13.8 reference phylogeny
758 (McDonald et al., 2012), using default parameters, to generate an inclusive phylogeny. An
759 observation table of per-sample counts of Deblur sub-operational taxonomic units were output
760 into BIOM format for analyses (n = 1511 samples). Outside of Qiita, we assigned taxonomy to
761 denoised reads using QIIME2's feature-classifier, `classify-sklearn`, using the GreenGenes 13.8
762 pre-fitted sklearn-based classifier (i.e., 99% OTUs, 515f/806r region of sequences), and default
763 parameters (Bokulich et al., 2018; Bolyen et al., 2019).

764

765 SourceTracker analyses

766 SourceTracker 2.0.1 (<http://github.com/biota/sourcetracker2>) was used to predict the proportion
767 of microbial source environment contributions to a sink using a Bayesian classification model
768 together with Gibbs sampling (Knights et al., 2011). The Deblur 150-bp observation table
769 consisting of 1511 food samples was used as the set of source environments and the Rheumatoid
770 Arthritis (RA) data set consisting of 49 fecal samples was used as the sink. All source and sink
771 samples were rarefied to 2000 sequences per sample before source-tracking and doubleton
772 ASVs were removed. Leave-one-out cross-validation was used to predict the source samples with
773 heterogeneity from all other food categories. After source sample filtering a total of 346 samples
774 representing a total of 25 broad food categories were retained. Food microbial source
775 contributions were then predicted for RA samples and the difference in food contribution before
776 and after diet intervention was calculated and compared by diet recommendations.

777 Metabolite Extraction

778 Homogenized samples were incubated for 40 min at -20°C and centrifuged (Eppendorf centrifuge
779 5418, Hamburg, Germany) at 20,000 rpm for 15 min at 4°C. 400 µL of supernatant were
780 transferred to a 96-well deep well plate and dried by centrifugal evaporation (Labconco Acid-
781 Resistant Centrivap Concentrator, Missouri, USA). Dried extracts were reconstituted in 150 µL of
782 resuspension solution (50% methanol with 2 µM sulfadimethoxine), then vortexed for 2 min and
783 sonicated for 5 min in a bath water (Branson 5510, Connecticut, USA). Resuspended extracts
784 were then centrifuged for 15 min at 20,000 rpm and 4°C (Thermo SORVALL LEGEND RT,
785 Germany) and transferred to a 96-well shallow well plate, and diluted either 5x or 10x to avoid
786 saturating the MS detector.

787 Liquid Chromatography - Mass Spectrometry

788 Food extracts were analyzed using a UltiMate 3000 UHPLC system (Thermo Scientific, Waltham,
789 Ma) equipped with a reverse phase C18 column, prepended with a guard cartridge (Kinetex, 100
790 x 2.1 mm, 1.7 µm particles size, 100 Å pore size; Phenomenex, Torrance, CA, USA), at a column
791 compartment temperature of 40°C. Samples were chromatographically separated with a constant
792 flow rate of 0.5 ml / min using the following gradient: 1.5 min isocratic at 5% B, up to 100% B in 8
793 min, 3 min isocratic at 100% B, back to 5% B in 0.5min and then 1.5min isocratic at 5% B (A: H₂O
794 + 0.1% formic acid; B: Acetonitrile (ACN) + 0.1% formic acid (LC-MS grade solvents, Fisher
795 Chemical, Hampton, United States)).

796 The UHPLC system was coupled to a Maxis Q-TOF Impact II mass spectrometer (Bruker
797 Daltonics, Bremen, Germany) equipped with an electrospray ionization source. MS spectra were
798 acquired in positive ionization mode using Data Dependent Acquisition (DDA) with a mass range
799 of m/z 50–1500. The instrument was externally calibrated two times per day to 1.0 ppm mass
800 accuracy using ESI-L Low Concentration Tuning Mix (Agilent Technologies, Waldbronn,
801 Germany). Hexakis (m/z 622.029509; (1H, 1H, 2H difluoroethoxy)phosphazene (Synquest
802 Laboratories, Alachua, FL)) was used for lock mass correction. MS/MS spectra were acquired for
803 the top 5 ions in each MS1 spectrum, with active exclusion after 2 spectra (maintained for 30
804 seconds). Known contaminants as well as lock mass values commonly used with this instrument

805 were added to an exclusion list (m/z values listed): 144.49–145.49; 621.00–624.10; 643.80–
806 646.00; 659.78–662.00; 921.0–925.00; 943.80–946.00; 959.80–962.00.

807 Raw high resolution mass spectrometry data files were converted to open source .mzXML
808 format using Bruker DataAnalysis software after lock mass correction (m/z 622.0290). Raw data
809 files as well as converted .mzXML files were uploaded to MassIVE (publicly available under
810 unique identifier MSV000084900) and further analyzed on Global Natural Product Social
811 Molecular Networking (GNPS) (<https://gnps.ucsd.edu>), as described below.

812 MS2 Data Processing

813 *FDR estimation*

814 False discovery rate (FDR) estimation was calculated using Passatutto analysis workflow in
815 GNPS (Scheubert et al. 2017; Wang et al. 2016). FDR estimation was used to determine the
816 cosine value required with a minimum of 5 matched peaks to achieve an FDR of 1%. See the
817 Data availability section for accession information.

818 *Molecular networking using GNPS:* Molecular networking analysis and library search were
819 performed using GNPS classical molecular networking release_18 (Wang et al. 2016). 3579
820 .mzXML data files (available at MassIVE ID MSV000084900) were included in the analysis. The
821 data were filtered by removing all MS/MS peaks within +/- 17 m/z of the precursor m/z . MS/MS
822 spectra were window filtered by choosing only the top 5 peaks in the +/- 50 m/z window throughout
823 the spectrum. The data was then clustered with MS-Cluster with a parent mass tolerance of 0.02
824 m/z and an MS/MS fragment ion tolerance of 0.02 m/z to create consensus spectra. Further,
825 consensus spectra that contained less than 2 spectra were discarded. A network was then
826 created where edges were filtered to have a cosine score above 0.65 (slight variation per study
827 based on FDR calculation) and more than 5 matched peaks. Further edges between two nodes
828 were kept in the network if and only if each of the nodes appeared in each other's respective top
829 10 most similar nodes. The spectra in the network were then searched against the GNPS spectral
830 libraries. The library spectra were filtered in the same manner as the input data. All matches kept
831 between network spectra and library spectra were required to have the same cosine score and
832 minimum matched peaks as for library search. Version release 18 was used to process all studies
833 with the exception of the Covid-19 dataset, which was processed with identical methods and
834 version 23.

835 Molecular networks were visualized in the GNPS browser as well as with the freely
836 available program Cytoscape (version 3.5.1) (Shannon et al., 2003).

837 *Interpreted spectral rate calculation*

838 The levels of interpretation are delineated as follows: A spectral match between an MS/MS
839 spectrum from human or food data with a library spectrum constitutes a *molecular ID* and
840 determines the initial percent of interpreted spectra, which is also equivalent to the annotation
841 rate of the dataset. A spectral match between MS/MS spectra in human and reference samples
842 (by performing molecular networking of the datasets together and identifying nodes with overlap
843 between the two groups) indicates a *potential source*. Matches between human and food data
844 therefore implicate food as the potential source of the molecule. Food reference data are referred
845 to in two main categories: the Global FoodOmics dataset (GFOP; broad range of foods and

846 beverages) and study specific food (SSF; foods and/or beverages known to be consumed by
847 some participants). The last level of interpretation is based on connectivity within a molecular
848 family, which allows us to infer *structural relatedness* or *possible metabolism* of food derived
849 compounds.

850 Food reference data and human data were organized into separate groups in the
851 molecular networking analysis. The annotation and interpreted spectral rates were calculated
852 using R (3.6.3) and the *tidyr* and *dplyr* packages. We first calculated percent annotation rate, or
853 molecular ID, for all studies (stool, plasma, etc.) (i.e. # of stool nodes with a molecular ID / total #
854 of stool nodes). Spectral matches between food reference data and human MS data (overlap
855 between the two groups) provides the next level of information, referred to as the interpreted
856 spectral rate (i.e. # of nodes found in food and stool data / total # of stool nodes), indicating a
857 potential food source.

858 For molecules without annotations to reference libraries, we wanted to measure the
859 potential to explain their presence using molecular networking. By removing single loops in each
860 dataset and comparing metabolites that shared a component index with an annotated compound,
861 we were able to identify molecules that belong to the same molecular family to infer their potential
862 classification, and calculate the interpreted spectral rate by dividing unannotated molecules that
863 network with annotated ones by total metabolites within each sample type. Overlap between
864 sample types was again assessed to understand contributions due to co-networking of molecules
865 across sample types, increasing our ability to explain unannotated molecules found in our
866 datasets. Visualizations were generated using *graphics* and *beeswarm* packages, and significant
867 differences were calculated using Welch's *t*-tests (*stats::t.test*), Welch's F-test
868 (*onewaytests::welch.test*), and Games-Howell (*rstatix::games_howell_test*) for multiple
869 comparisons, as appropriate, with multiple comparisons correction using Tukey's method. All data
870 are expressed as the mean \pm standard error and considered significant if $P < 0.05$ unless
871 otherwise stated.

872 For example, for GNPS molecular networking analyses test datasets were consistently
873 placed in group 1 (G1) (and G2 for paired datasets, such as stool and plasma) and Global
874 FoodOmics data were placed in group 4 (G4). SSF were consistently placed in G3 when used.
875 The common nodes between G1 and G4 represent the overlap and potential enhancement of
876 information, directly from the reference dataset. The improvement is thus measured by the
877 difference in the overlap of G1 and G4 divided by the total nodes in G1 versus the # of annotations
878 in G1 divided by the total nodes in G1. The "propagation" refers to the counting of nodes within
879 connected components in molecular families which capture three types of additional information:
880 1) unannotated compounds found only in G1 that network with an annotated compound found in
881 G4 (could be an annotated molecule observed only in G4 or in G4 and G1), 2) unannotated
882 compounds found only in G1, but in the same molecular family with an unannotated food
883 compound (G4), or 3) unannotated compounds found only in G1, but in the same molecular family
884 with an annotated food compound (G4). The increase shown for Total is taking into account the
885 # of unique nodes from the three different types of molecular connectivity. The second is the
886 largest contributor.

887 Metadata inference - food count generation

888 Food counts were calculated as the number of consensus nodes in the molecular networking
889 results that match to food samples. Consensus nodes were required to match to all of the relevant
890 experiment groups (sample type, GFOP, optionally SSF) and not match to any of the other
891 experiment groups. All source file names corresponding to the filtered consensus nodes were
892 matched to the GFOP file names and metadata to derive counts of the foods at different levels of
893 the food hierarchy. Infrequent food types that occurred less often than water (presumed blank)
894 were removed to filter out sporadic random matches.

895 For the flow diagrams the food counts for the complete datasets were calculated at
896 different levels of the metadata hierarchy. Flow diagrams were generated in Python (version 3.8)
897 using Pandas (version 0.25.3) (McKinney, 2010), NumPy (version 1.18.1) (van der Walt et al.,
898 2011), and floweaver (version 2.0.0a5) (Lupton and Allwood, 2017).

899 Diet validation with RA dataset

900 The food counts at the fifth hierarchy level were extracted for each individual raw file and used to
901 construct a feature table. The occurrences were summed across groups (diet intervention),
902 divided by the total number of samples in each group, respectively and the difference was
903 calculated. These differences were then compared with the ITIS diet recommendations by food
904 category. Foods were grouped into one of four categories: avoid, include, restricted, and not
905 specified.

906 Diet diary entries were tabulated across over 200 categories and the closest matches to
907 the food categories identified by MS were identified. The corresponding diet data was tabulated
908 based on the number of times a category was reported during the three time points prior to the
909 diet change (pre-intervention) and the three time points prior to sample collection of the final time
910 point (post-intervention; during the intervention). The sum of the three days per diet category was
911 then divided by the total number of samples in the pre vs. post sample group, respectively (to
912 account for missing self-reported information). Three days were chosen as a representation of
913 foods most likely to be detected (Johnson et al., 2019). Categories were matched as closely as
914 possible to those in the FoodOmics ontology.

915 Dataset descriptions

916 All human datasets were processed by LC-MS/MS on high resolution mass spectrometers, in
917 positive ionization mode.

918 Data were collected for the following studies using a QTOF mass spectrometer and similar
919 methods as those outlined above: American Gut (MSV000081981), Children with Medical
920 Complexity (MSV000084610), Rotarix vaccine response (MSV000084218), Malawi legume
921 supplement (MSV000081486), IBD_1 (MSV000082431), IBD_individual (MSV000079115),
922 Fermented food consumption (MSV000081171) (Taylor et al., 2020). The Sleep deprivation
923 (MSV000083759; IRB 15-0282), centenarian (MSV000084591; IRB 180478), and Legume
924 supplementation (MSV000084663) studies were analyzed using the methods described above
925 and described in (Gauglitz et al., 2020a). The LP Infant (MSV000083462; MSV000083463),
926 IBD_seed (MSV000082221), IBD_biobank (MSV000079777), IBD_2 (MSV000084775), IBD_200

927 (MSV000084908), IBD_biopsy (MSV000082220), Gout (MSV000084908), Adult Saliva
928 (MSV000083049) datasets were collected as described previously (Gauglitz et al., 2020b).

929 The datasets for the impact of diet on RA (MSV000084556) and Alzheimer's disease
930 (MSV000085256) were collected with similar methods on a Q-exactive Orbitrap mass
931 spectrometer (Thermo Scientific). The Alzheimer samples include Alzheimer's Disease and
932 elderly controls, and were drawn in the early morning after fasting for at least 6 hours.

933 The food and plasma data for the Covid-19 study (MSV000085505; MSV000085537) were
934 collected at the University of São Paulo, Brazil, as described below: Plasma samples were
935 collected from patients with laboratory confirmed Covid-19 who were admitted to the Special Unit
936 for the Treatment of Infectious Diseases (UETDI) at the General Hospital of the Medical School
937 of Ribeirão Preto (HC-FMRP-USP). Previously, clarifications to patients occurred both orally and
938 in writing, based on the printed text of the Free and Informed Consent Form, which contained the
939 general proposal of the study, the procedures for obtaining the samples, the risks and benefits.
940 In addition, they were assured about confidentiality of their name, personal data and the possibility
941 of giving up their participation at any time. Following the signature, patients received a copy of
942 the informed consent form. The following were included: 1) Patients diagnosed with Covid-19 in
943 moderate, severe or critical forms and in need of hospital treatment; 2) Over 18 years old; 3) At
944 least 50 kg of body weight; 4) Admission electrocardiogram without changes in rhythm and with
945 QT interval <450 ms; 5) normal serum levels of Ca²⁺ and K⁺; 6) If a woman, between 18 and 50
946 years old, negative β-HCG test on admission. Patients were excluded who: 1) have the mild forms
947 of SARS-CoV-2; 2) pregnant; 3) unable to understand the information contained in the Free and
948 Informed Consent Form (ICF).

949 Sample preparation: For the Covid-19 plasma samples, aliquots of 20 μL were transferred
950 to eppendorf tubes and 120 μL of cold extracting solution, MeOH: MeCN (1: 1, v/v) was added.
951 After orbital shaking for 1 min (Gehaka AV-2 Shaker, São Paulo, Brazil), the samples were left at
952 -20 °C for 30 minutes and then centrifuged for 10 min at 20000 × g at 4 °C (Centrifuge Boeco
953 Germany M-240R, Germany). An aliquot of the organic phase (120 μL) was transferred to another
954 eppendorf tube and evaporated to dryness in a rotary vacuum concentrator for 60 min, at 30 °C
955 (Analytica, Christ RVC2-18, São Paulo). The residues were resuspended in 80 μL of H₂O and
956 centrifuged (10 min, 5000 ×g, 4 °C), an aliquot of 5 μL was injected.

957 Mass spectrometry data collection plasma sample extracts were chromatographically
958 separated with anHPLC (Shimadzu, Tokyo, Japan), coupled with a micrOTOF-Q II mass
959 spectrometer (Bruker Daltonics, Boston, MA, USA) equipped with an ESI source and a
960 quadrupole-time of flight analyzer (qTOF, Bruker Daltonics Inc., Billerica, MA, USA). For
961 chromatographic analyses, we employed a Kinetex C18 column (1.7 μm, 100 × 2.1 mm)
962 (Phenomenex, Torrance, CA, USA) kept at 40 °C, with a flow rate of 0.3 mL/min. A linear gradient
963 was applied: 0-1.5 min isocratic at 5% B, 1.5-9.5 min 100% B, 9.5-12 min isocratic at 100% B,
964 12-12.5 min 5% B, 12.5-14 min 5% B; where mobile phase A is water with 0.1% formic acid (v/v)
965 and phase B is acetonitrile 0.1% formic acid (v/v) (LC-MS grade solvents). The MS data were
966 acquired in positive mode using an MS range of *m/z* 50–1500. The equipment was calibrated with
967 trifluoroacetic acid (TFA) every day, and internally during each run. The MS parameters were
968 established as follows: end plate offset, 450 V; capillary voltage, 3500 V; nebulizer gas pressure,
969 4.0 Bar; dry gas flow, 9 L/min; dry temperature, 220 °C.

970 For data dependent acquisition the five most abundant ions per MS1 scan were
971 fragmented and the spectra collected. MS/MS active exclusion was set after 2 spectra and
972 released after 30 seconds. A fragmentation exclusion list was set: m/z 144.49-145.49; 621.00-
973 624.10; 643.80- 646.00; 659.78-662.00; 921.0-925.00; 943.80-946.00; 959.80-962.00 to exclude
974 known contaminants and infused lock mass compounds. A process blank was run every 5
975 samples; 5 μL of a standard mix [Paclitaxel 1 mg L^{-1} , and Diazepam 1 mg L^{-1}] (Sigma-Aldrich,
976 Saint Louis, Missouri, US) in 50% MeOH (LC-MS grade solvents) was injected every 5 samples.
977 All MS data were analyzed with Bruker Compass DataAnalysis 4.3 software (Bruker Daltonics,
978 Boston, MA, USA).

979 A metadata file was created grouping all available clinical information from patients with
980 laboratory confirmed Covid-19 and essential analysis specifications. The MS/MS data were
981 calibrated with an internal standard (TFA), converted to mzXML files using MSConvert from the
982 ProteoWizard software (Chambers et al., 2012) and then uploaded into the Global Natural
983 Products Social Molecular Networking web-platform (<https://gnps.ucsd.edu/>). All MS data
984 (.mzXML files) and metadata (.txt file) are publically available via GNPS/MassIVE
985 (<https://massive.ucsd.edu/>) under accession number MSV000085373.

986

987 Supplement

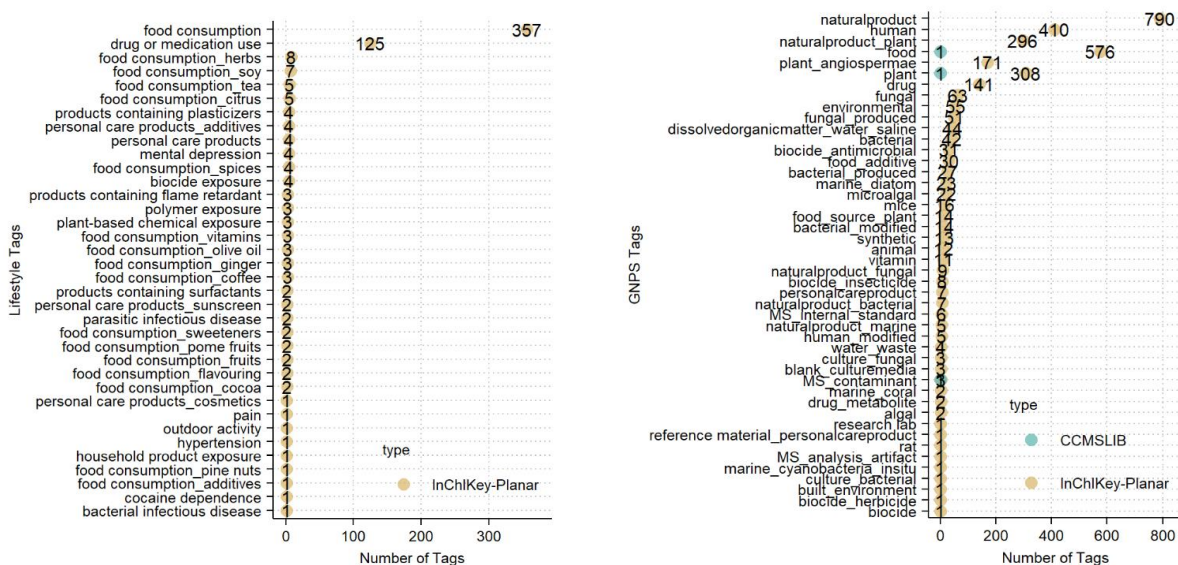
988 **Table S1. Metadata** [available also on MassIVE under ID MSV000084900]

989 **Table S2.** Overview of public studies used in analysis. Each sample type represents an individual
990 dataset.

Study	Sample Type	SSF	Num samples	Massive ID
GFOP3500	Food	N/A	3527	MSV000084900
Sleep study	Fecal; Plasma	yes (197)	98 (F); 371 (P)	MSV000083759
Centenarian	Fecal; Plasma	yes (38)	91 (F); 50 (P)	MSV000084591
Impact of diet on RA	Fecal; Plasma	yes (12)	51 (F); 60 (P)	MSV000084556
LP Infant	Fecal; Oral; Skin	yes (58)	492(F); 461(O); 461(S)	MSV000083462; MSV000083463
Children with Medical Complexity	Fecal	yes (24)	95	MSV000084610
American Gut	Fecal		2123	MSV000081981
Fermented food consumption	Fecal		276	MSV000081171
Malawi legume supplement	Fecal	yes (14)	1131	MSV000081486
Rotarix vaccine response	Fecal		118	MSV000084218
IBD_1	Fecal		40	MSV000082431
IBD_individual	Fecal		5	MSV000079115
IBD_seed	Fecal		334	MSV000082221
IBD_biobank	Fecal		95	MSV000079777
IBD_2	Fecal		206	MSV000084775
IBD_200	Fecal		203	MSV000084908
Alzheimer's disease	Plasma; CSF		78 (P); 116 (CSF)	MSV000085256
Covid-19 Brazil	Plasma	yes (60)	46	MSV000085505; MSV000085537
IBD_biopsy	Tissue		135	MSV000082220
Gout	Serum		39	MSV000084908
Adult saliva	Saliva		89	MSV000083049

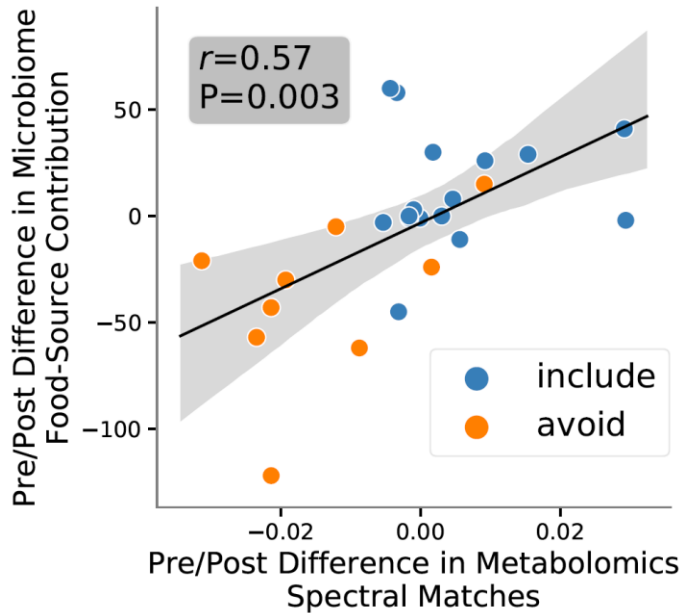
Legume supplementation	Urine	yes (15)	5	MSV000084663
-------------------------------	-------	----------	---	--------------

991
992



993
994
995
996
997
998
999
1000
1001
1002
1003

Figure S1. GNPS tag and GNPS Lifestyle Tag distribution for the Global FoodOmics reference data set (GNPS task ID: f1a1f3a61aca416a9b3687d72488da7f). Annotated MS/MS spectra were assigned planar InChIKeys, and at least one tag. Spectra can be assigned multiple tags, indicating multiple potential sources. 1131 total unique planar InChIKeys with at least one GNPS tag. **a.** Lifestyle tags and **b.** GNPS tags.



1004
1005 **Figure S2.** Linear regression scatter plot of difference in food contributions for metabolite spectral
1006 match (x-axis) and microbes by source tracking prediction (y-axis) before vs. after diet intervention
1007 compared by diet recommendation of avoid (orange) or include (blue). Correlation evaluated by
1008 Pearson correlation coefficient.

1009 References

- 1010 Amir, A., McDonald, D., Navas-Molina, J.A., Kopylova, E., Morton, J.T., Zech Xu, Z., Kightley,
1011 E.P., Thompson, L.R., Hyde, E.R., Gonzalez, A., et al. (2017). Deblur Rapidly Resolves Single-
1012 Nucleotide Community Sequence Patterns. *mSystems* 2.
- 1013 Aron, A.T., Gentry, E.C., McPhail, K.L., Nothias, L.-F., Nothias-Esposito, M., Bouslimani, A.,
1014 Petras, D., Gauglitz, J.M., Sikora, N., Vargas, F., et al. (2020). Reproducible molecular
1015 networking of untargeted mass spectrometry data using GNPS. *Nat. Protoc.* 15, 1954–1991.
- 1016 Barabási, A.-L., Menichetti, G., and Loscalzo, J. (2020). The unmapped chemical complexity of
1017 our diet. *Nature Food* 1, 33–37.
- 1018 Bokulich, N.A., Kaehler, B.D., Rideout, J.R., Dillon, M., Bolyen, E., Knight, R., Huttley, G.A., and
1019 Gregory Caporaso, J. (2018). Optimizing taxonomic classification of marker-gene amplicon
1020 sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome* 6.
- 1021 Bolyen, E., Rideout, J.R., Dillon, M.R., Bokulich, N.A., Abnet, C.C., Al-Ghalith, G.A., Alexander,
1022 H., Alm, E.J., Arumugam, M., Asnicar, F., et al. (2019). Reproducible, interactive, scalable and
1023 extensible microbiome data science using QIIME 2. *Nat. Biotechnol.* 37, 852–857.
- 1024 Bono, H. (2020). All of gene expression (AOE): An integrated index for public gene expression
1025 databases. *PLoS One* 15, e0227076.
- 1026 Bouslimani, A., Melnik, A.V., Xu, Z., Amir, A., da Silva, R.R., Wang, M., Bandeira, N.,
1027 Alexandrov, T., Knight, R., and Dorrestein, P.C. (2016). Lifestyle chemistries from phones for

- 1028 individual profiling. *Proc. Natl. Acad. Sci. U. S. A.* *113*, E7645–E7654.
- 1029 Bustamante, M.F., Agustín-Perez, M., Cedola, F., Coras, R., Narasimhan, R., Golshan, S., and
1030 Guma, M. (2020). Design of an anti-inflammatory diet (ITIS diet) for patients with rheumatoid
1031 arthritis. *Contemp Clin Trials Commun* *17*, 100524.
- 1032 Center for Food Safety, and Nutrition, A. (2019). *Food Defect Levels Handbook*.
- 1033 Chambers, M.C., Maclean, B., Burke, R., Amodei, D., Ruderman, D.L., Neumann, S., Gatto, L.,
1034 Fischer, B., Pratt, B., Egertson, J., et al. (2012). A cross-platform toolkit for mass spectrometry
1035 and proteomics. *Nat. Biotechnol.* *30*, 918–920.
- 1036 Gauglitz, J.M., Aceves, C.M., Aksenov, A.A., Aleti, G., Almaliti, J., Bouslimani, A., Brown, E.A.,
1037 Campeau, A., Caraballo-Rodríguez, A.M., Chaar, R., et al. (2020a). Untargeted mass
1038 spectrometry-based metabolomics approach unveils molecular changes in raw and processed
1039 foods and beverages. *Food Chem.* *302*, 125290.
- 1040 Gauglitz, J.M., Morton, J.T., Tripathi, A., Hansen, S., Gaffney, M., Carpenter, C., Weldon, K.C.,
1041 Shah, R., Parampil, A., Fidgett, A.L., et al. (2020b). Metabolome-Informed Microbiome Analysis
1042 Refines Metadata Classifications and Reveals Unexpected Medication Transfer in Captive
1043 Cheetahs. *mSystems* *5*.
- 1044 Gonzalez, A., Navas-Molina, J.A., Kosciolk, T., McDonald, D., Vázquez-Baeza, Y.,
1045 Ackermann, G., DeReus, J., Janssen, S., Swafford, A.D., Orchanian, S.B., et al. (2018). Qiita:
1046 rapid, web-enabled microbiome meta-analysis. *Nat. Methods* *15*, 796–798.
- 1047 Haug, K., Cochrane, K., Nainala, V.C., Williams, M., Chang, J., Jayaseelan, K.V., and
1048 O'Donovan, C. (2020). MetaboLights: a resource evolving in response to the needs of its
1049 scientific community. *Nucleic Acids Res.* *48*, D440–D444.
- 1050 Horai, H., Arita, M., Kanaya, S., Nihei, Y., Ikeda, T., Suwa, K., Ojima, Y., Tanaka, K., Tanaka,
1051 S., Aoshima, K., et al. (2010). MassBank: a public repository for sharing mass spectral data for
1052 life sciences. *J. Mass Spectrom.* *45*, 703–714.
- 1053 Huang, R., Zhu, H., Shinn, P., Ngan, D., Ye, L., Thakur, A., Grewal, G., Zhao, T., Southall, N.,
1054 Hall, M.D., et al. (2019). The NCATS Pharmaceutical Collection: a 10-year update. *Drug Discov.*
1055 *Today* *24*, 2341–2349.
- 1056 Janssen, S., McDonald, D., Gonzalez, A., Navas-Molina, J.A., Jiang, L., Xu, Z.Z., Winker, K.,
1057 Kado, D.M., Orwoll, E., Manary, M., et al. (2018). Phylogenetic Placement of Exact Amplicon
1058 Sequences Improves Associations with Clinical Information. *mSystems* *3*.
- 1059 Jarmusch, A.K., Wang, M., Aceves, C.M., Advani, R.S., Aguire, S., Aksenov, A.A., Aleti, G.,
1060 Aron, A.T., Bauermeister, A., Bolleddu, S., et al. (2019). Repository-scale Co- and Re-analysis
1061 of Tandem Mass Spectrometry Data.
- 1062 Johnson, A.J., Vangay, P., Al-Ghalith, G.A., Hillmann, B.M., Ward, T.L., Shields-Cutler, R.R.,
1063 Kim, A.D., Shmagel, A.K., Syed, A.N., Personalized Microbiome Class Students, et al. (2019).
1064 Daily Sampling Reveals Personalized Diet-Microbiome Associations in Humans. *Cell Host*
1065 *Microbe* *25*, 789–802.e5.
- 1066 Knights, D., Kuczynski, J., Charlson, E.S., Zaneveld, J., Mozer, M.C., Collman, R.G., Bushman,

- 1067 F.D., Knight, R., and Kelley, S.T. (2011). Bayesian community-wide culture-independent
1068 microbial source tracking. *Nat. Methods* 8, 761–763.
- 1069 Kyle, J.E., Crowell, K.L., Casey, C.P., Fujimoto, G.M., Kim, S., Dautel, S.E., Smith, R.D., Payne,
1070 S.H., and Metz, T.O. (2017). LIQUID: an open source software for identifying lipids in LC-
1071 MS/MS-based lipidomics data. *Bioinformatics* 33, 1744–1746.
- 1072 Lupton, R.C., and Allwood, J.M. (2017). Hybrid Sankey diagrams: Visual analysis of
1073 multidimensional data for understanding resource use. *Resour. Conserv. Recycl.* 124, 141–151.
- 1074 Martin, C.L., Montville, J.B., Steinfeldt, L.C., Omolewa-Tomobi, G., Heendeniya, K.Y., Adler,
1075 M.E., and Moshfegh, A.J. (2012). USDA Food and Nutrient Database for Dietary Studies 2011--
1076 2012: Documentation and User Guide. Beltsville, MD: US Department of Agriculture.
1077 Agricultural Research Service, USDA Food Surveys Research Group.
- 1078 McDonald, D., Price, M.N., Goodrich, J., Nawrocki, E.P., DeSantis, T.Z., Probst, A., Andersen,
1079 G.L., Knight, R., and Hugenholtz, P. (2012). An improved Greengenes taxonomy with explicit
1080 ranks for ecological and evolutionary analyses of bacteria and archaea. *ISME J.* 6, 610–618.
- 1081 McDonald, D., Hyde, E., Debelius, J.W., Morton, J.T., Gonzalez, A., Ackermann, G., Aksenov,
1082 A.A., Behsaz, B., Brennan, C., Chen, Y., et al. (2018). American Gut: an Open Platform for
1083 Citizen Science Microbiome Research. *mSystems* 3.
- 1084 McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Proceedings of the
1085 9th Python in Science Conference, (SciPy), pp. 56–61.
- 1086 Ono, H., Ogasawara, O., Okubo, K., and Bono, H. (2017). RefEx, a reference gene expression
1087 dataset as a web tool for the functional analysis of genes. *Sci Data* 4, 170105.
- 1088 Quinn, R.A., Nothias, L.-F., Vining, O., Meehan, M., Esquenazi, E., and Dorrestein, P.C. (2017).
1089 Molecular Networking As a Drug Discovery, Drug Metabolism, and Precision Medicine Strategy.
1090 *Trends Pharmacol. Sci.* 38, 143–154.
- 1091 Sawada, Y., Nakabayashi, R., Yamada, Y., Suzuki, M., Sato, M., Sakata, A., Akiyama, K.,
1092 Sakurai, T., Matsuda, F., Aoki, T., et al. (2012). RIKEN tandem mass spectral database
1093 (ReSpect) for phytochemicals: a plant-specific MS/MS-based data resource and database.
1094 *Phytochemistry* 82, 38–45.
- 1095 Scheubert, K., Hufsky, F., Petras, D., Wang, M., Nothias, L.-F., Dührkop, K., Bandeira, N.,
1096 Dorrestein, P.C., and Böcker, S. (2017). Significance estimation for large scale metabolomics
1097 annotations by spectral matching. *Nat. Commun.* 8, 1494.
- 1098 Schmid, R., Petras, D., Nothias, L.-F., Wang, M., Aron, A.T., Jagels, A., Tsugawa, H., Rainer,
1099 J., Garcia-Aloy, M., Dührkop, K., et al. (2020). Ion Identity Molecular Networking in the GNPS
1100 Environment.
- 1101 Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N.,
1102 Schwikowski, B., and Ideker, T. (2003). Cytoscape: a software environment for integrated
1103 models of biomolecular interaction networks. *Genome Res.* 13, 2498–2504.
- 1104 Sicherer, S.H., and Sampson, H.A. (2006). 9. Food allergy. *J. Allergy Clin. Immunol.* 117, S470–
1105 S475.

- 1106 Song, S.J., Amir, A., Metcalf, J.L., Amato, K.R., Xu, Z.Z., Humphrey, G., and Knight, R. (2016).
1107 Preservation Methods Differ in Fecal Microbiome Stability, Affecting Suitability for Field Studies.
1108 *mSystems* 1.
- 1109 Sprecher, K.E., Ritchie, H.K., Burke, T.M., Depner, C.M., Smits, A.N., Dorrestein, P.C.,
1110 Fleshner, M., Knight, R., Lowry, C.A., Turek, F.W., et al. (2019). Trait-like vulnerability of higher-
1111 order cognition and ability to maintain wakefulness during combined sleep restriction and
1112 circadian misalignment. *Sleep* 42.
- 1113 Sumner, L.W., Amberg, A., Barrett, D., Beale, M.H., Beger, R., Daykin, C.A., Fan, T.W.-M.,
1114 Fiehn, O., Goodacre, R., Griffin, J.L., et al. (2007). Proposed minimum reporting standards for
1115 chemical analysis Chemical Analysis Working Group (CAWG) Metabolomics Standards
1116 Initiative (MSI). *Metabolomics* 3, 211–221.
- 1117 Taylor, B.C., Lejzerowicz, F., Poirel, M., Shaffer, J.P., Jiang, L., Aksenov, A., Litwin, N.,
1118 Humphrey, G., Martino, C., Miller-Montgomery, S., et al. (2020). Consumption of Fermented
1119 Foods Is Associated with Systematic Differences in the Gut Microbiome and Metabolome.
1120 *mSystems* 5.
- 1121 Thompson, L.R., Sanders, J.G., McDonald, D., Amir, A., Ladau, J., Locey, K.J., Prill, R.J.,
1122 Tripathi, A., Gibbons, S.M., Ackermann, G., et al. (2017). A communal catalogue reveals Earth's
1123 multiscale microbial diversity. *Nature* 551, 457–463.
- 1124 Vinaixa, M., Schymanski, E.L., Neumann, S., Navarro, M., Salek, R.M., and Yanes, O. (2016).
1125 Mass spectral databases for LC/MS- and GC/MS-based metabolomics: State of the field and
1126 future prospects. *Trends Analyt. Chem.* 78, 23–35.
- 1127 van der Walt, S., Colbert, S.C., and Varoquaux, G. (2011). The NumPy Array: A Structure for
1128 Efficient Numerical Computation. *Computing in Science Engineering* 13, 22–30.
- 1129 Wang, M., Carver, J.J., Phelan, V.V., Sanchez, L.M., Garg, N., Peng, Y., Nguyen, D.D.,
1130 Watrous, J., Kapon, C.A., Luzzatto-Knaan, T., et al. (2016). Sharing and community curation of
1131 mass spectrometry data with Global Natural Products Social Molecular Networking. *Nat.*
1132 *Biotechnol.* 34, 828–837.
- 1133 Watrous, J., Roach, P., Alexandrov, T., Heath, B.S., Yang, J.Y., Kersten, R.D., van der Voort,
1134 M., Pogliano, K., Gross, H., Raaijmakers, J.M., et al. (2012). Mass spectral molecular
1135 networking of living microbial colonies. *Proc. Natl. Acad. Sci. U. S. A.* 109, E1743–E1752.
- 1136 Wishart, D.S., Feunang, Y.D., Marcu, A., Guo, A.C., Liang, K., Vázquez-Fresno, R., Sajed, T.,
1137 Johnson, D., Li, C., Karu, N., et al. (2018). HMDB 4.0: the human metabolome database for
1138 2018. *Nucleic Acids Res.* 46, D608–D617.
- 1139