1    American Gut: an Open Platform for Citizen-Science Microbiome Research

2

3    Daniel McDonald[a,%], Embriette Hyde[a,%], Justine W. Debelius[a], James T. Morton[a], Antonio

4    Gonzalez[a], Gail Ackermann[a], Alexander A. Aksenov[b,c], Bahar Behsaz[d], Caitriona Brennan[a],

5    Yingfeng Chen[e],  Lindsay DeRight Goldasich[a], Pieter C. Dorrestein[b,c], Robert R. Dunn[f], Ashkaan

6    K. Fahimipour[g], James Gaffney[a], Jack A Gilbert[h,i,j,k], Grant Gogul[a], Jessica L. Green[g], Philip

7    Hugenholtz[l], Greg Humphrey[a], Curtis Huttenhower[m,n], Matthew A. Jackson[o], Stefan Janssen[a],

8    Dilip V. Jeste[p,q], Lingjing Jiang[a], Scott T. Kelley[14], Dan Knights[r,s], Tomasz Kosciolek[a], Joshua

9    Ladau[t], Jeff Leach[u], Clarisse Marotz[a], Dmitry Meleshko[v], Alexey V. Melnik[b,c], Jessica L.

10   Metcalf[w], Hosein Mohimani[x], Emmanuel Montassier[r,y], Jose Navas-Molina[a], Tanya T.

11   Nguyen[p,q], Shyamal Peddada[z], Pavel Pevzner[b,d,aa], Katherine S. Pollard[t], Gholamali

12   Rahnavard[m,n], Adam Robbins-Pianka[bb], Naseer Sangwan[j], Joshua Shorenstein[a], Larry Smarr[d,aa,cc],

13   Se Jin Song[a], Timothy Spector[o], Austin D. Swafford[aa], Varykina G. Thackray[dd], Luke R.

14   Thompson[ee], Anupriya Tripathi[a], Yoshiki Vazquez-Baeza[a], Alison Vrbanac[a], Paul

15   Wischmeyer[ff,gg], Elaine Wolfe[a], Qiyun Zhu[a], The American Gut Consortium, Rob Knight[a,d,aa,#]

16

17   [a] Department of Pediatrics, University of California San Diego, La Jolla, CA, USA

18   [b] Collaborative Mass Spectrometry Innovation Center, University of California, San Diego, La

19   Jolla, CA, USA

20   [c] Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego,

21   La Jolla, CA, USA

22   [d] Department of Computer Science and Engineering, University of California, San Diego, La

23   Jolla, CA, USA

24    e Department of Biology, San Diego State University, San Diego, CA, USA

25    f Department of Applied Ecology, North Carolina State University, Raleigh, NC, USA

26    g Biology and the Built Environment Center; University of Oregon, Eugene, OR, USA

27    h Department of Surgery, University of Chicago, Chicago, IL, USA

28    i Institute for Genomic and Systems Biology, University of Chicago, Chicago, IL, USA

29    j Department of Biosciences, Argonne National Laboratory, Chicago, IL, USA

30    k Marine Biology Laboratory, University of Chicago, Chicago, IL, USA

31    l Australian Centre for Ecogenomics, School of Chemistry and Molecular Biosciences, The

32    University of Queensland, Brisbane, QLD, Australia

33    m Harvard T. H. Chan School of Public Health, Boston, MA, USA

34    n The Broad Institute of MIT and Harvard, Boston, MA, USA

35    o Department of Twin Research and Genetic Epidemiology, King's College London, London,

36    UK

37    p Departments of Psychiatry and Neurosciences, University of California San Diego, La Jolla,

38    CA, USA

39    q Sam and Rose Stein Institute for Research on Aging, and Center for Healthy Aging, University

40    of California San Diego, La Jolla, CA, USA

41    r Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN,

42    USA

43    s Biotechnology Institute, University of Minnesota, Minneapolis, MN, USA

44    t The Gladstone Institutes, University of California, San Francisco, CA, USA

45    u Human Food Project, Terlingua, TX, USA

1

46    ᵛ St. Petersburg State University, Center for Algorithmic Biotechnology, Saint Petersburg,

47    Russia

48    ʷ Department of Animal Science, Colorado State University, Fort Collins, CO, USA

49    ˣ Department of Computational Biology, Carnegie Mellon University, Pittsburgh, PA, USA

50    ʸ Université de Nantes, Microbiotas Hosts Antibiotics and Bacterial Resistances (MiHAR),

51    Nantes, France

52    ᶻ Department of Biostatistics, University of Pittsburgh, Pittsburgh, PA, USA.

53    ᵃᵃ Center for Microbiome Innovation, University of California, San Diego, La Jolla, CA, USA

54    ᵇᵇ Department of Computer Science, University of Colorado Boulder, Boulder, Colorado, USA

55    ᶜᶜ California Institute for Telecommunications and Information Technology (Calit2), University

56    of California San Diego, La Jolla, CA, USA

57    ᵈᵈ Department of Reproductive Medicine, University of California San Diego, La Jolla, CA,

58    USA

59    ᵉᵉ Southwest Fisheries Science Center, National Oceanic and Atmospheric Administration, La

60    Jolla, CA, USA

61    ᶠᶠ Department of Anesthesiology and Surgery, Duke University School of Medicine, Durham,

62    NC, USA

63    ᵍᵍ Duke Clinical Research Institute, Duke University School of Medicine, Durham, NC, USA

64

65    Running Title (max 54c): American Gut: an Open Platform for Microbiome Research

66    # Address correspondence to Rob Knight, robknight@ucsd.edu

67    % D.M. and E.H. contributed equally to this work.

68    Abstract word count: 203

69     Main text word count: 3280

70

71     **Abstract**: Although much work has linked the human microbiome to specific phenotypes and

72     lifestyle variables, data from different projects have been challenging to integrate and the extent

73     of microbial and molecular diversity in human stool remains unknown. Using standardized

74     protocols from the Earth Microbiome Project and sample contributions from over 10,000 citizen-

75     scientists, together with an open research network, we compare human microbiome specimens

76     primarily from the USA, UK, and Australia to one another and to environmental samples. Our

77     results show an unexpected range of beta-diversity in human stool microbiomes as compared to

78     environmental samples, demonstrate the utility of procedures for removing the effects of

79     overgrowth during room-temperature shipping for revealing phenotype correlations, uncover

80     new molecules and kinds of molecular communities in the human stool metabolome, and

81     examine emergent associations among the microbiome, metabolome, and the diversity of plants

82     that are consumed (rather than relying on reductive categorical variables such as veganism,

83     which have little or no explanatory power). We also demonstrate the utility of the living data

84     resource and cross-cohort comparison to confirm existing associations between the microbiome

85     and psychiatric illness, and to reveal the extent of microbiome change within one individual

86     during surgery, providing a paradigm for open microbiome research and education.

87

88     **Importance:** We show that a citizen-science, self-selected cohort shipping samples through the

89     mail at room temperature recaptures many known microbiome results from clinically collected

90     cohorts and reveals new ones. Of particular interest is integrating n=1 study data with the

91     population data, showing that the extent of microbiome change after events such as surgery can

92    exceed differences between distinct environmental biomes, and the effect of diverse plants in the

93    diet which we confirm with untargeted metabolomics on hundreds of samples.

94

95    **Introduction**

96         The human microbiome plays a fundamental role in human health and disease. While

97    many studies link microbiome composition to phenotypes, we lack understanding of the

98    boundaries of bacterial diversity within the human population, and the relative importance of

99    lifestyle, health conditions, and diet, to underpin precision medicine or to educate the broader

100   community about this key aspect of human health.

101        We launched the American Gut Project (AGP; http://americangut.org) in November of

102   2012 as a collaboration between the Earth Microbiome Project (EMP) (1) and the Human Food

103   Project (HFP; http://humanfoodproject.com/) to discover the kinds of microbes and microbiomes

104   "in the wild" via a self-selected citizen-scientist cohort. The EMP is tasked with characterizing

105   the global microbial taxonomic and functional diversity, and the HFP is focused on

106   understanding microbial diversity across human populations. As of May 2017, the AGP included

107   microbial sequence data from 15,096 samples from 11,336 human participants, totaling over 467

108   million (48,599 unique) 16S rRNA V4 gene fragments ("16S"). Our project informs citizen-

109   scientist participants about their own microbiomes by providing a standard report (fig 1A) and

110   resources to support human microbiome research, including an online course (Gut Check:

111   Exploring Your Microbiome; *https://www.coursera.org/learn/microbiome*). AGP deposits all de-

112   identified data into the public domain on an ongoing basis without access restrictions (table S1).

113   This reference database characterizes the diversity of the industrialized human gut microbiome

114   on an unprecedented scale, reveals novel relationships with health, lifestyle, and dietary factors,

115   and establishes the AGP resource and infrastructure as a living platform for discovery.

116

117    **Results**

118    **Cohort characteristics.** AGP participants primarily reside in the United States (n=7,860).

119    However, interest in the AGP rapidly expanded beyond the US to United Kingdom (n=2,518),

120    and Australia (*n*=321), with 42 other countries or territories also represented (fig 1A; table S1).

121    Participants in the US inhabit urban (*n*=7,317), rural (*n*=29), and mixed (*n*=98) communities

122    (2010 US Census data based on participant zip codes), and span greater ranges of age, race, and

123    ethnicity than other large-scale microbiome projects (2–6). Because the AGP is crowdsourced

124    and self-selected, and subjects generally support the cost of sample processing, the population is

125    unrepresentative in several important respects, including having lower prevalence of smoking

126    and obesity, higher education and income (fig S1A), and underrepresentation of Hispanic and

127    African American communities (table S1); generalization of the results is cautioned. Targeted

128    and population-based studies will be crucial for filling these cohort gaps (Supplemental text).

129         Using a survey modified from (7, 8), participants reported general health status, disease

130    history, and lifestyle data (table S2, supplemental text). In accordance with our IRB, all survey

131    questions were optional (median per-question response 70.9%; table S2). Additionally, 14.8% of

132    participants completed a validated picture-based food frequency questionnaire (FFQ)

133    (VioScreen; *http://www.viocare.com/vioscreen.html*), and responses correlated well with primary

134    survey diet responses (table S2).

135         We sought to minimize errors and misclassifications well-known to occur in self-reported

136    data (9). Survey responses relied on controlled vocabularies. For analyses, we trimmed numeric

137    entries at extremes (e.g., weight over 200kg or below 2.5kg) and excluded obviously incorrect

138    answers (e.g., infants drinking alcohol) and samples for which necessary data were not supplied

139    (e.g., missing zip code data for spatial analyses); see supplement for details. We focused our

140    primary investigative efforts on a "healthy adult" subset ($n$=3,942) of individuals aged 20-69

141    with BMIs ranging between 18.5–30 kg/m$^2$, no self-reported history of inflammatory bowel

142    disease, diabetes, or antibiotic use in the past year, and at least 1,250 16S sequences/sample (fig

143    1B, S1B).

144            The two largest populations in the dataset (US and UK) differed significantly in alpha-

145    diversity, with Faith's phylogenetic diversity (PD) higher in UK samples (*13*) (Mann Whitney

146    $p<1\text{x}10^{-15}$; fig 1C). One balance (10) (a log-ratio compositional transform) explained most of the

147    taxonomic separation between US and UK samples (AUC=77.7% ANOVA $p=1.01\text{x}10^{-78}$,

148    $F$=386.85) (fig S1C, table S3). To understand how these two populations differed from others,

149    we compared adult AGP samples (predominantly from industrialized regions) to samples from

150    adults living traditional lifestyles (6, 11, 12). As previously observed (6), samples from industrial

151    and traditional populations separated in Principal Coordinates Analysis (PCoA) space of

152    unweighted UniFrac distances (13) (fig S1D). They show greater variation within industrial

153    populations than within traditional populations (2) and facile separation based on microbial

154    taxonomy (industrial vs. non-industrial agrarian: AUC=98.9%, ANOVA $p=1.52\text{x}10^{-260}$,

155    $F$=1265.8; industrial vs. hunter-gatherer: AUC=99.5%, ANOVA $p=4.48\text{x}10^{-227}$, $F$=1092.35) (fig

156    1D, table S3).

157

158    **Removal of bacterial blooms.** An important practical question is whether self-collected

159    microbiome samples can match those from better-controlled studies. Most AGP samples are

160    stools collected on dry swabs and shipped without preservative to minimize costs and avoid

161    exposure to toxic preservatives. *E. coli* and a few other taxa grow in transit, so based on data

162    from controlled storage studies as previously described (14) we removed sOTUs (sub-OTUs

163    (15); median of 7.9% of sequences removed per sample) shown to bloom.

164          We further characterized the impact of these organisms through culturing, HPLC-MS

165    analysis of cultured isolates, and shotgun metagenomics of the primary samples and storage

166    controls (16). Culturing primary specimens stored at -80°C (US: $n=116$; UK: $n=73$; other: $n=25$)

167    showed a strong correlation between the fraction of sequences reported as blooms in 16S

168    sequencing and positive microbial growth following overnight incubation in aerobic conditions

169    (fig 2A). Culture supernatants were characterized using HPLC-MS; most metabolites in these

170    supernatants were absent from the primary specimens (fig 2B, C, method details in SI). We

171    sequenced draft genomes of 169 isolates; of these, 65 contained the exact *E. coli* 16S sequence in

172    the published bloom filter (14). To characterize the impact of the 16S bloom filter, we computed

173    effect sizes over the participant covariates and technical parameters for 9,511 individual

174    participant samples, including and excluding blooms (complete list table S2; comparisons to (17,

175    18) in supplementary text), and observed tight correlations for both unweighted (fig 2D, Pearson

176    $r=0.91$, $p=3.76 \times 10^{-57}$; Spearman $r=0.90$, $p=9.45 \times 10^{-55}$) and weighted UniFrac (fig 2E, Pearson

177    $r=0.42$, $p=1.71 \times 10^{-6}$; Spearman $r=0.58$, $p=1.03 \times 10^{-9}$). An outlier on the quantitative metric

178    (weighted UniFrac) is present and corresponds to a variable representing the fraction of bloom

179    reads in a sample.

180

181    **Novel taxa and microbiome configurations.** To understand human microbiome diversity, we

182    placed AGP samples in the context of the EMP (1). Building on earlier work revealing a striking

183    difference between host-associated and environmental microbiomes (19), we found that the

184     diversity of microbiomes associated with the human gut (just one vertebrate) occupies a vast

185     extent of the microbiome diversity of the planet (fig 3A).

186         Inserting the sOTU fragments of AGP and EMP samples into a Greengenes (20)

187     reference phylogenetic tree using SEPP (21) (fig 3B) showed that the AGP population harbored

188     much broader microbial diversity than the Human Microbiome Project (5). Both datasets are

189     dwarfed by the breadth of bacterial and archaeal phylogenetic diversity in environmental

190     samples. Examining sOTUs over increasing numbers of samples, we observed a reduction in the

191     discovery rate of novel sOTUs starting around 3,000 samples, emphasizing the need for focused

192     sampling efforts outside the present AGP population (fig 3C). The importance of sample size for

193     detecting novel microbes and microbiomes is apparent when contrasted against Yatsunenko et al.

194     (6), which contained hundreds of samples from three distinct human populations at ~1 million

195     sequences/sample (fig 3D). This effect is magnified in beta-diversity analysis, where the AGP

196     has saturated the configuration space, and new samples are not "distant" from existing samples

197     (fig 3E). To encourage community engagement with sOTUs found in the AGP, we adapted the

198     EMP "trading cards" for sOTUs (figs 3F, S2).

199

200     **Temporal and spatial analyses.** Longitudinal samples are required for understanding human

201     microbiome dynamics (22). We examined 565 individuals who contributed multiple samples and

202     observed an increasing trend of intrapersonal divergence with time. Still, over time individuals

203     resemble themselves more than others, even after one year (fig 4A).

204         We tested whether patterns in individual longitudinal sample sets could be better

205     explained when placed in the context of the AGP by integrating samples collected from: a) a

206     time series of 58 time points from one subject (described as "LS"), prior to and following a large

207     bowel resection, b) 2 time points from 121 patients in an intensive care unit (ICU) (23),  c)

8

208    samples from the "extreme" diet study from David et al. (24), and d) samples from the Hadza

209    hunter-gatherers for additional context (25). Through the longitudinal sampling of LS, dramatic

210    pre- and post-microbial configuration changes that exceeded the span of microbial diversity

211    associated with the AGP population were observed (fig 4E, animated in (26)). After surgery,

212    subject's samples more closely resembled those of ICU patients (Kruskal Wallis H=79.774,

213    p=4.197x$^{-19}$, fig S2A-C), and showed a persistent state change upon return to the AGP fecal

214    space. Remarkably, the UniFrac distance between the samples immediately prior to and

215    following the surgery was almost identical to the distance between a marine sediment sample and

216    a plant rhizosphere sample (unweighted UniFrac distance of 0.78). Furthermore, the observed

217    state change in LS is not systematically observed in the extreme diet study (fig S2D;

218    PERMANOVA n.s. when controlling for individual). Despite extensive dietary shifts, these

219    subjects do not deviate from the background AGP context.

220         Recent reports suggest that the microbes of bodies (8), like those of homes (27), are

221    influenced mostly by local phenomena rather than regional biogeography (28), and accordingly

222    we observed only weak geographic associations with sOTUs (fig 4B), no significant distance-

223    decay relationships (fig 4C), and, with Bray-Curtis distance, only a weak effect at neighborhood

224    sizes of ca. 100km (Mantel $r$=0.036, Benjamini-Hochberg adjusted $p$=0.03) to 1,000km (Mantel

225    $r$=0.016, Benjamini-Hochberg adjusted $p$=0.03).

226

227    **Dietary plant diversity.** The self-reported dietary data suggested, unexpectedly, that the number

228    of unique plant species a subject consumes is associated microbial diversity, rather than self-

229    reported categories such as "vegan" or "omnivore" (fig 2D, E). Principal Components Analysis

230    of FFQ responses (fig 5A) revealed clusters associated with diet types such as "vegan."

9

231    However, these dietary clusters did not significantly relate to microbiome configurations (fig 5B;

232    Procrustes fig 5A, $M^2$=0.988). We therefore characterized the impact of dietary plant diversity on

233    the microbial community.

234        Using balances (10), we identified several putative short-chain fatty acid (SCFA)

235    fermenters associated with eating more than 30 types of plants, including sOTUs putatively of

236    the species *Faecalibacterium prausnitzii* and of the genus *Oscillospira* (29) (AUC=68.5%,

237    ANOVA *p*=8.9x10$^{-39}$, *F*=177.2) (fig 5E, table S3). These data suggest community-level changes

238    associated with microbial fermentation of undigested plant components. Because bacteria differ

239    in their carbohydrate-binding modules and enzymes that hydrolyze diverse substrates in the gut

240    (30), a diet containing various types of dietary fibers and resistant starches likely supports a more

241    diverse microbial community (31, 32).

242        To test these effects in the stool metabolome, we performed HPLC-MS annotation and

243    annotation propagation (33, 34) on a subset of fecal samples (*n*=219) preferentially selecting

244    individuals at the extremes of plant type consumption, i.e. eating <10 or >30 different types of

245    plants per week. Several fecal metabolites differed between the two groups, with one key

246    discriminating feature annotated as octadecadienoic acid (annotation level 2 according to the

247    2007 metabolomics initiative, (35)). Further investigation using authentic standards revealed that

248    the detected feature was comprised of multiple isomers, including linoleic acid (LA) and

249    conjugated linoleic acid (CLA). CLA abundance was significantly higher in individuals

250    consuming > 30 types of plants, and those consuming more fruits and vegetables generally, (fig

251    5D, 1-sided *t*-test; p < 10$^{-5}$), but did not correlate with dietary CLA consumption as determined

252    by the FFQ (dietary fig 5C; Spearman *r* < 0.16; *p* > 0.15). CLA is a known end-product of LA

253    conversion by lactic acid bacteria in the gut, such as *Lactobacillus plantarum* (36) and

254   *Bifidobacterium* spp. (37). FFQ-based dietary levels of LA and MS-detected LA did not differ

255   significantly between groups (fig S3), suggesting that their different microbiomes may

256   differentially convert LA to CLA. Several other putative octadecadienoic acid isomers were also

257   detected (fig 5F), some strongly correlated with plant consumption. Determining these

258   compounds' identities as well as their origin and function may uncover new links between the

259   diet, microbiome, and health.

260

261   **Molecular novelty in the human gut metabolome.** Our untargeted HPLC-MS approach

262   allowed us to search for novel molecules in the human stool metabolome, parallel to our search

263   for novelty in microbes and microbiome configurations described above. Bacterial N-acyl

264   amides were recently shown to regulate host metabolism by interacting with G-protein-coupled

265   receptors (GPCRs) in the murine gastrointestinal tract, mimicking host-derived signaling

266   molecules (38). These agonistic molecules regulate metabolic hormones and glucose

267   homeostasis as efficiently as host ligands. Manipulating microbial genes that encode metabolites

268   eliciting host cellular responses could enable new drugs or treatment strategies for many major

269   diseases, including diabetes, obesity, and Alzheimer's disease: roughly 34% of all marketed

270   drugs target GPCRs (39). We observed N-acyl amide molecules previously hypothesized but

271   unproven to be present in the gut (38) (fig 6, S4), as well as new N-acyl amides (fig 6).

272       Levels of two N-acyl amides, annotated as commendamide (*m/z* 330.2635, fig S4B) and

273   N-3-OH-palmitoyl ornithine (*m/z* 387.3220, fig S4C), positively correlated with a self-reported

274   medical diagnosis of thyroid disease (Kruskal–Wallis, FDR $p$=0.032, $p$=2.48x10$^{-3}$, $\chi$2=11.99; N-

275   3-OH-palmitoyl ornithine; Kruskal–Wallis, FDR $p$=0.048, $p$=5.63x10$^{-3}$, $\chi$2=10.35). Conversely,

276   glycodeoxycholic acid (m/z 450.3187) was significantly higher in individuals not reporting

277   thyroid disease diagnosis (Kruskal–Wallis; FDR $p$=1.28x10$^{-4}$, $p$= 4.41x10$^{-7}$, $\chi$2=29.27). This

11

278  cholic acid is produced through microbial dehydroxylation, again linking gut microbiota to

279  endocrine function (40, 41).

280  Finally, we compared metabolome diversity to 16S diversity in the samples selected for

281  dietary plant diversity and a second set of samples selected to explore antibiotic effects (*n*=256

282  individuals who self-reported not having taken antibiotics in the past year (*n*=117), or having

283  taken antibiotics in the past month (*n=*139); participants were matched for age, BMI, and

284  country). By computing a collector's curve of observed molecular features in both cohorts (fig

285  6K, 6M), we observe that, paradoxically, individuals who had taken antibiotics in the past month

286  (*n*=139) had significantly greater molecular diversity (Kruskal Wallis, $H$=255.240, $p$=1.87x10$^{-57}$)

287  than those who had not taken antibiotics in the past year (*n*=117), and differed in molecular beta-

288  diversity (fig 6K inset), suggesting that antibiotics promote unique metabolomes that result from

289  differing chemical and microbial environments in the gut. Notably, the diversity relationships of

290  this set are not reflected in 16S diversity (fig 6L, 6N), where antibiotic use shows decreased

291  diversity (Kruskal Wallis $H$=3983.839, $p$=0.0). Within the dietary plant diversity cohort, we

292  observed a significant increase (Kruskal Wallis, $H$=897.106, $p$=4.17x10$^{-197}$) in molecular alpha

293  diversity associated with a high diversity of plant consumption (*n*=42) compared to low plant

294  diversity (*n*=43), a relationship also observed in 16S diversity, where high dietary plant diversity

295  increased 16S alpha diversity (Kruskal Wallis, $H$=65.817, $p$=4.947x$^{-16}$).

296

297  **A living dataset.** The AGP is dynamic, with samples arriving from around the world daily. This

298  allows a living analysis, similar to continuous molecular identification and annotation revision in

299  the Global Natural Products Molecular Networking (GNPS) database (34). Although the analysis

300  presented here represents a single snapshot, samples continued to arrive during manuscript

301  preparation. For example, after we defined the core "healthy" sample set, an exploratory analysis

12

302    using matched controls was performed by collaborators to test for correlations between mental

303    illness and microbiome composition (as reported in (42, 43)). By analyzing mental illness status

304    (depression, schizophrenia, post-traumatic stress disorder (PTSD) and bipolar disorder – four of

305    the most disabling illnesses per World Health Organization (44)) reported by AGP participants

306    ($n$=125) against matched 1:1 healthy controls ($n$=125), we observed a significant partitioning

307    using PERMANOVA in weighted UniFrac ($p$=0.05, $pseudo$-$F$=2.36). These findings were

308    reproducible within US residents ($n$=122, p=0.05, $pseudo$-$F$=2.58), UK residents ($n$=112,

309    $p$=0.05, $pseudo$-$F$=2.16), women ($n$=152, $p$=0.04, $pseudo$-$F$=2.35), and people 45 years of age or

310    younger ($n$=122, $p$=0.05, $pseudo$-$F$=2.45). We also reproduce some previously reported

311    differentially abundant taxa in Chinese populations using our UK subset (42, 45)(table S3). This

312    shows that multi-cohort replication is possible within the AGP (additional detail supplemental

313    text).

314

315    **Discussion**

316        The AGP provides an example of a successful crowdfunded citizen science project that

317    facilitates human microbiome hypothesis generation and testing on an unprecedented scale,

318    provides a free data resource derived from over 10,000 human-associated microbial samples, and

319    both recaptures known microbiome results and yields new ones. Ongoing living data efforts,

320    such as the AGP, will allow researchers to document and potentially mitigate the effects of a

321    slow but steady global homogenization driven by increased travel, lifespans, and access to

322    similar diets and therapies, including antibiotics. Because the AGP is a subproject of the EMP

323    (1), all samples were processed using the publicly available and widely used EMP protocols to

324    facilitate meta-analyses, as highlighted above. Further example applications include assessing the

325    stability of AGP runs over time, comparing the AGP population to fecal samples collected from

326    a fecal transplant study (46) and an infant microbiome time series (47), the latter using different

327    DNA sequencing technology, to highlight how this context can provide insight (48).

328         A unique aspect of the AGP is the open community process of assembling the Research

329    Network and analyzing these data, which are released immediately on data generation. Analysis

330    details are shared through a public forum (GitHub, https://github.com/knightlab-

331    analyses/american-gut-analyses). Scientific contributions to the project were made through a

332    geographically diverse Research Network represented herein as the American Gut Consortium,

333    established prior to project launch and which has grown over time. This model allows a "living

334    analysis" approach, embracing new researchers and analytical tools on an ongoing basis (e.g.,

335    Qiita (*Web:http://qiita.microbio.me*) and GNPS (34)). Examples of users of the AGP as a

336    research platform include educators at several universities, UC San Diego Athletics, and the

337    American Gastroenterological Association (AGA). Details on projects using the AGP

338    infrastructure can be found in the supplement.

339         To promote public data engagement, we aimed to broaden the citizen science experience

340    obtained by participating in AGP by "gamifying" the data and separately by developing an

341    online forum for microbiome data discussion and discovery. The gamification introduces

342    concepts of beta-diversity and challenges users to identify clusters of data in principal

343    coordinates space (http://csb.cs.mcgill.ca/colonyb/). The forum, called Gut Instinct

344    (http://gutinstinct.ucsd.edu), enables participants to share lifestyle-based insights with one

345    another. Participants also have the option to share their AGP sample barcodes, which will help us

346    uncover novel contextual knowledge. Gut Instinct now has over 1,050 participants who have

347    collectively created over 250 questions. Participants will soon design and run their own

348     investigations using controlled experiments to further understand their own lifestyle and the AGP

349     data.

350         The AGP therefore represents a unique citizen-science dataset and resource, providing a

351     rich characterization of microbiome and metabolome diversity at the population level. We

352     believe the community process for involving participants from sample collection through data

353     analysis and deposition will be adopted by many projects harnessing the power of citizen science

354     to understand the world around and within our own bodies.

355

356     **Materials and methods**

357     **Participant Recruitment and Sample Processing.** Participants signed up for the project

358     through Indiegogo (https://www.indiegogo.com/) and later, FundRazr (http://fundrazr.com/). A

359     contribution to the project was made to help offset the cost of sample processing and sequencing

360     (typically $99 per sample; no requirement to contribute if another party was covering the

361     contribution). All participants were consented under an approved Institutional Review Board

362     human research subjects protocol, either from the University of Colorado Boulder (protocol #12-

363     0582; December 2012 - March 2015) or the University of California, San Diego (protocol

364     #141853; February 2015 - present). The IRB-approved protocol specifically allows for public

365     deposition of all data that is not personally identifying and for return of results to participants

366     (fig. 1A).

367

368     Self-reported metadata were collected through a web portal

369     (http://www.microbio.me/americangut). Samples were collected using BBL Culture Swabs

370     (Becton, Dickinson and Company; Sparks, MD) and returned by mail. Samples were processed

371   using the EMP protocols. Briefly, the V4 region of the 16S rRNA gene was amplified with

372   barcoded primers and sequenced as previously described (49). Sequencing prior to August 2014

373   was done using the 515f/806r primer pair with the barcode on the reverse primer (50);

374   subsequent rounds were sequenced with the updated 515f/806rB primer pair with the barcode on

375   the forward read (51). Sequencing batches 1-19 and 23-49 were sequenced using an Illumina

376   MiSeq; sequencing for 20 and 21 were performed with an Illumina HiSeq Rapid Run and round

377   22 was sequenced with an Illumina HiSeq High Output.

378

379   **16S Data Processing.** The 16S sequence data were processed using a sequence variant method,

380   Deblur v1.0.2 (52) trimming to 125nt (otherwise default parameters), to maximize the specificity

381   of 16S data; a trim of 125nt was used because one sequencing round in the American Gut used

382   125 cycles while the rest used 150. Following processing by Deblur, previously recognized

383   bloom sequences were removed (14). The Deblur sub Operational Taxonomic Units (sOTUs)

384   were inserted into the Greengenes 13_8 (53) 99% reference tree using SEPP (54). Taxonomy

385   was assigned using an implementation of the RDP classifier (55) as implemented in QIIME2

386   (56). Multiple rarefactions were computed, with the minimum being 1250 sequences per sample

387   with the analyses using the 1250 set except where noted explicitly. Diversity calculations were

388   computed using scikit-bio 0.5.1 with the exception of UniFrac (57) which was computed using

389   an unpublished algorithmic variant, Striped UniFrac (*https://github.com/biocore/unifrac*), which

390   scales to larger datasets and produces identical results to previously published UniFrac

391   algorithms.

392

16

393  **Metadata Curation.** To address the self-reported nature of the AGP data and ongoing nature of

394  the project, basic filtering was performed on the age, height, weight, and body mass index

395  (BMI). Height and weight were gated to only consider heights between 48 cm and 210 cm, and

396  weight between 2.5 kg and 200 kg. BMI calculations using values outside this range were not

397  considered. We assumed age was misreported by any individual who reported a birth date after

398  their sample was collected. We also assumed age was misreported for participants who reported

399  an age of less than 4 years, but height over 105 cm, weight over 20 kg, or any alcohol

400  consumption. Values assumed to be incorrect were dropped from analyses (fig S1B).

401

402  **Sample Selection**. Analyses in the manuscript were performed on a subset of the total AGP

403  samples. A single fecal sample was selected for each participant with at least one fecal sample

404  that amplified to 1250 sequences per sample unless otherwise noted. Priority was given to

405  samples that were associated with VioScreen (*http://www.viocare.com/vioscreen.html)* metadata.

406

407  The samples used for analysis and subsets used in various analyses are described in table S2.

408  Briefly, we defined the healthy subset (*n*=3,942) as adults aged 20-69 years with a BMI between

409  18.5 and 30 kg/m$^2$ who reported no history of inflammatory bowel disease or diabetes and no

410  antibiotic use in the last year. There were 1,762 participants who provided results for the

411  VioScreen Food Frequency Questionnaire (FFQ; *http://www.viocare.com/vioscreen.html*). The

412  meta-analysis with non-Western samples (*n*=4,643) included children over the age of 3, adults

413  with a BMI of between 18.5 and 30 kg/m$^2$, and no reported history of inflammatory bowel

414  disease, diabetes, or antibiotic use in the last year.

415

416     **Population Level Comparisons.** Population level comparisons were calculated for all American

417     Gut participants living in the United States. BMI categorization was only considered for adults

418     over the age of twenty, since the description of BMI in children is based on their age and sex.

419     Education level was considered for adults over the age of 25. This threshold was used to match

420     the available data from the US Census Bureau

421     (*https://www.census.gov/content/dam/Census/library/publications/2016/demo/p20- 578.pdf*). The

422     percentage of the American Gut participants was calculated as the fraction of individuals who

423     reported results for that variable. US population data is from the 2010 census

424     (*https://www.census.gov/prod/cen2010/briefs/c2010br-03.pdf*), US Census bureau reports

425     (*https://www.census.gov/content/dam/Census/library/publications/2016/demo/p20- 578.pdf*),

426     Centers for Disease Control reports on obesity

427     (*https://www.cdc.gov/nchs/data/hus/2015/058.pdf*), diabetes (57, 58), IBD

428     (*http://www.cdc.gov/ibd/ibd-epidemiology.htm*), smoking

429     (*https://www.cdc.gov/tobacco/data_statistics/fact_sheets/adult_data/cig_smoking/index.ht*

430     *m*), and a report from the Williams Institute (*http://williamsinstitute.law.ucla.edu/wp-*

431     *content/uploads/How-Many-Adults-Identify-as-Transgender-in-the-United-States.pdf*) (table S2).

432

433     **Within American Gut Alpha- and Beta-Diversity Analyses.** OTU tables generated in the

434     primary processing step were rarefied to 1,250 sequences per sample. Shannon, Observed OTU,

435     and PD whole tree diversity metrics were calculated as the mean of ten rarefactions using QIIME

436     (56, 59). Alpha-diversity for single metadata categories was compared with a Kruskal-Wallis

437     test. Unweighted UniFrac distance between samples was tested with PERMANOVA (60) and

438     permuted *t*-tests in QIIME.

439

440 **Balances.** The goal of this analysis was to design two-way classifiers to classify samples and

441 sOTUs. This will allow us to identify sOTUs that are strongly associated with a given

442 environment. To do this while accounting for issues due to compositionality, we used balances

443 (61) constructed from Partial Least Squares (62).

444

445 First the sOTU table was centered log-ratio (CLR) transformed with a pseudocount of 1. Partial

446 least squares discriminant analysis (PLS-DA) was then performed on this sOTU table using a

447 single PLS component, using a binary categorical variable as the response and the CLR

448 transformed sOTU table as the predictor. This PLS component represented an axis, which

449 assigns scores to each OTU according to how strongly associated they are to each class. An

450 sOTU with a strong negative score indicates an association for the one category, which we will

451 denote as the negative category. An sOTU with a strong positive score indicates that sOTU is

452 strongly associated with the other category, which we will denote as the positive category.

453

454 We assumed that PLS scores associated with each OTU were normally distributed. Specifically

455

456 $score(x_{pos}^{(i)}) \sim N(\mu_{pos}, \sigma_{pos}^2)$

457 $score(x_{neg}^{(i)}) \sim N(\mu_{neg}, \sigma_{neg}^2)$

458 $score(x_{null}^{(i)}) \sim N(\mu_{null}, \sigma_{null}^2)$

459

460 Where $\mu_{null} \approx 0$, $\mu_{neg} < 0$ and $\mu_{pos} > 0$. To obtain estimates of these normal distributions,

461 Gaussian Mixture Models with three Gaussians were fitted from the PLS scores. Thresholds

462 were determined from the intersection of Gaussians. The OTUs with PLS scores less than the

19

463    intersection $N(\mu_{null}, \sigma_{null}^2)$ and $N(\mu_{neg}, \sigma_{neg}^2)$ are classified to be associated with the negative

464    category.   The OTUs with PLS scores greater than the intersection $N(\mu_{null}, \sigma_{null}^2)$ and

465    $N(\mu_{pos}, \sigma_{pos}^2)$ are classified to be associated with the positive category.

466    The balance was constructed as follows

467

468    $$b = \sqrt{\frac{|x_{pos}||x_{neg}|}{|x_{pos}| + |x_{neg}|}} \, log\left(\frac{g(x_{pos})}{g(x_{neg})}\right)$$

469

470    From this balance, we calculated receiver operator characteristic (ROC) curves and AUC to

471    assess the classification accuracy, and ran ANOVA to assess the statistical significance. The

472    dimensionality was shrunk through some initial filtering (an sOTU must have at least 50 reads,

473    must exist in at least 20 samples except where noted, and have a variance over 10 to remove

474    sOTUs that do not appear to change), so that the number of samples is greater than the number of

475    sOTUs to reduce the likelihood of over-fitting. This technique was used to investigate

476    differences due to plant consumption, country of residence and western vs non-western and was

477    consistently applied with the exception that a filter of 5 samples was used for the western vs.

478    non-western analysis due to group sample sizes.

479

480    Balances on plant consumption were constructed using Partial Least Squares.  Only samples

481    from people who consumed less than 10 types of plants a week or more than 30 types of plants a

482    week were considered.

483

484    **Meta-analysis of samples from the American Gut and from individuals living agrarian and**

485    **hunter-gatherer lifestyles.** A meta-analysis compared fecal samples collected from healthy

20

486     individuals that were 3 years of age or older and included in the AGP data set to a previously

487     published 16S rRNA V4 region data set that included healthy people living an industrialized,

488     remote agrarian or hunter-gatherer lifestyle (63–65). The AGP subset of healthy individuals was

489     determined by filtering by the metadata columns "subset_antibiotic", "subset_ibd",

490     "subset_diabetes", and for individuals over the age of 16 years "subset_bmi". All datasets were

491     processed using the Deblur pipeline as noted above, with the exception that all reads in the meta-

492     analysis, including AGP data, were trimmed to 100nt to accommodate the read length in

493     Yatsunenko et al (63). Bloom reads as described above were removed from all samples. We used

494     Striped UniFrac as noted above to estimate beta-diversity (unweighted UniFrac) and EMPeror

495     software (66) version 0.9 to visualize principal coordinates. We used a non-parametric

496     PERMANOVA with 999 permutations to test for significant differences in fecal microbiomes

497     associated with industrialized, remote agrarian, and hunter-gatherer lifestyles. All AGP samples

498     were considered to be from people living an industrialized lifestyle. Balances were constructed

499     from Partial Least Squares to assess the differences between the hunter-gather vs. industrialized

500     populations and the remote farmers vs industrialized populations.

501

502     **Spatial Autocorrelation.** We sought to investigate distance-decay patterns – the relationship

503     between microbial community similarity and spatial proximity – among American Gut

504     participants, to determine the extent to which geographical distances could explain variation in

505     microbial community taxonomic compositions between participant pairs. The correlation

506     between community-level Bray-Curtis (67) distances and participants' spatial proximities (i.e.,

507     great-circle distances, km) was assessed using a Mantel test (68) with 1000 matrix permutations.

508     Analyses were conducted using the subset of participants located in the continental United States

509    that had not received antibiotics in the last year. Different neighborhood sizes were investigated

510    in order to detect the relevant spatial scale on which significant distance-decay patterns in

511    microbial community compositions emerged. To accomplish this, we computed distance-decay

512    relationships for a series of model adjacencies corresponding to neighborhood radiuses of 50,

513    100, 500, 1000, 2500, and 4500 km among participants, and adjusted $p$-values for multiple

514    comparisons using the Benjamini-Hochberg procedure (69). We also studied spatial correlations

515    in phylogenetic community dissimilarities, calculated as weighted normalized UniFrac distances,

516    using the procedure described above. Analyses were conducted in R statistical programming

517    environment.

518

519    The spatial autocorrelation of each individual taxon was assessed using Moran's $I$ statistic (70).

520    Taxa present in less than 10 samples were filtered, since these would not be sufficiently

521    powered. Analyses were conducted using binary spatial weight matrices, with neighborhoods of

522    $0 - 50$ km, $50 - 100$ km, and $100 - 250$ km. The different neighborhoods were useful for

523    detecting spatial autocorrelation at different scales. All spatial weights matrices were row-

524    standardized. We checked for spatial autocorrelation at three taxonomic ranks: class, genus, and

525    OTU. We also considered whether there was autocorrelation within subsets of individuals who

526    were under 20 years old and between 20 and 70 years old; those having IBD, no IBD, diabetes,

527    and no diabetes; and those who had taken antibiotics within the past week, year, or not within the

528    past year. The results presented above did not qualitatively depend on the subset of individuals

529    considered. Statistical significance was assessed using permutation tests, which were

530    implemented using a Markov Chain Monte Carlo algorithm. To assess each $p$-value, 100 chains

531    were run each starting from a different random permutation. Each chain had 1000 iterations. We

532    used Bonferroni corrections to correct for multiple comparisons, with an overall significance

533    level set to 0.05. Analyses were run using custom Java code, optimized for running many spatial

534    autocorrelation analyses on large data sets (71).

535

536    **Metadata cross-correlation**. To account for covariance among metadata for effect size and

537    variation analyses, we examined the correlation between individual metadata variables including

538    technical parameters. Groups in ordinal variables were combined if there were insufficient

539    sample size (e.g. people who reported sleeping less than 5 hours were combined with those who

540    reported sleeping 5 to 6 hours into a variable described as "Less than 6"). The same

541    transformations were used for effect size analysis. Any group with less than 25 total observations

542    was ignored during analysis; if this resulted in a metadata column having no groups, the column

543    was removed from analysis. The relationship between continuous and ordinal covariates was

544    calculated using Pearson's correlation. Ordinal and categorical covariates were compared using a

545    modified Cramer's V statistic (72). Continuous and categorical covariates were compared with a

546    Welch's T test (73). We treated used *1-R* as a distance between the covariates. Traversing the

547    resulting binary, weighted cluster tree starting at tip level into the direction of the root, i.e.

548    bottom-up, we grouped tips together that are members of the same subtree after covering a

549    distance of approximately 0.5 (branch length 0.29). A representative variable from each cluster

550    was selected for analysis (table S2).

551

552    **Effect Size Calculations**. Effect size was calculated on 179 covariates (including technical

553    parameters), selected from the cross-correlation (table S2). Ordinal groups with small sample

554    sizes at the extreme were collapsed as noted above. Individuals who reported self-diagnosis or

23

555    diagnosis from alternative practitioner for medical conditions were excluded from the analysis.

556    Any metadata variable with less than 50 observations per group or that made up less than 3% of

557    the total number of respondents was also excluded from the effect size analysis. Continuous

558    covariates were categorized into quartiles. For each one of the 179 variables, we applied the

559    mdFDR (74) methodology to test for the significance of each pairwise comparison among the

560    groups.  For each significant pairwise comparison, we computed the effect size using Cohen's $d$

561    (75), or the absolute difference between the mean of each group divided by the pooled standard

562    deviation. For analysis of diversity, we used Faith's Phylogenetic Diversity (alpha-diversity) and

563    weighted and unweighted UniFrac distances (beta-diversity).

564

565    **Variation analysis.** Using the methodology reported in the supplemental material for (76), we

566    computed Adonis (77) using 1000 permutations, over the sample sets used in the effect size

567    calculations as noted above, and applied Benjamini-Hochberg correction (FDR<0.1) to assess

568    drivers of variation in beta-diversity.

569

570    **Meta-analysis movie.** American Gut samples from all body sites were combined with data from

571    an infant time series (78), a fecal transplant study (79), and recent work characterizing the

572    microbiome of patients in the intensive care unit (80). The combination of the datasets in movie

573    S2 required that all sequences were trimmed to an even length of 125 nucleotides. All projects

574    except for the infant time series were sequenced using an Illumina instrument. In order to

575    combine the data, we expressed the Illumina and non-Illumina data through a common reference

576    database. Specifically, the Deblur sOTUs from the Illumina data were mapped against the

577    Greengenes (53) database (13_8 release) using 99% similarity; the associations between the

578    input sOTUs, and their cluster memberships, were used to construct an OTU table based on the

579    original sOTU per sample sequences counts (i.e., summing the counts for all sOTUs in a

580    common OTU). The infant time series data were picked using a closed reference OTU picking

581    approach against the same reference at the same similarity. The infant time series dataset

582    followed a closed reference OTU picking approach using 99% similarity. The resulting two

583    tables (from Illumina-generated data and the ITS dataset) were merged and analyzed using the

584    Greengenes 99% tree. The table was rarefied to 1,250 sequences per sample. Principal

585    coordinates projections were calculated based on unweighted UniFrac distance (57). The

586    principal coordinates analysis was visualized and animated in EMPeror 1.0.0-beta8-dev (66, 81).

587    The movie was captured in QuickTime (Apple, Cupertino, CA), and edited with Premiere Pro

588    (Adobe, San Jose, CA).

589

590    **Integration with the Earth Microbiome Project.** A precomputed 100nt Deblur BIOM table

591    representing the data in (82) was obtained from

592    (ftp://ftp.microbio.me/emp/release1/otu_tables/deblur/). 100nt Deblur tables were also obtained

593    from Qiita for Hadza fecal samples (Qiita study ID 11358, (83)), ICU microbiome samples (Qiita

594    study ID 2136, (80)), and a longitudinal series which includes samples immediately prior to and

595    following a large bowel resection (Qiita study ID 10283, EBI accession ERP105968,

596    unpublished); all samples were processed using the EMP Illumina 16S V4 protocol. The EMP

597    dataset used a minimum sOTU count of 25; the same threshold was applied to the other datasets

598    included prior to merge. Blooms as identified by (84) were removed from all samples. This

599    collection of BIOM tables was then merged yielding an OTU table representing 40,600 samples.

600    sOTUs were restricted to those already present in the EMP 100nt fragment insertion tree, which

25

601    represents 329,712 sOTUs. The table was then rarefied to 1000 sequences per sample, and

602    unweighted UniFrac computed using 768 processors with the aforementioned Striped algorithm.

603    Visualizations and animations were performed using EMPeror v1.0.0b12.dev0.

604

605    **Extreme diet study state assessment.** The sequence data from (85) were processed by Deblur to

606    assess 16S sOTUs in common with the AGP processing above. In order to assess a state

607    difference with PERMANOVA, we needed to control for sample independence within the

608    longitudinal sampling. To do so, we randomly selected one sample from each individual per diet,

609    computed PERMANOVA, and repeated the process 100 times. None of the trials produced a $p$-

610    value below 0.05.

611

612    **Vioscreen PCA and diet type Procrustes analysis.** Before performing Principal Component

613    Analysis (PCA) on the informal diet questions, Vioscreen variables that are categorical or

614    receive less than 90% response among the 1762 participants were excluded leaving 1596

615    participants. PCA was then performed using the Vioscreen information from these participants'

616    responses over 207 Vioscreen questions, and then colored by their types of diet as answered in

617    the AGP informal food survey. The coordinates from the PCA were extracted. For the same

618    samples, PCoA of unweighted UniFrac distances was computed on the 16S data subset from the

619    primary processing set. The coordinates from the PCA and the PCoA were assayed for a measure

620    of fitness using Procrustes as implemented in QIIME v1.9.1.

621

622    **Beta-diversity added.** To assess added beta-diversity, we applied the technique used in (86)

623    figure 3. Specifically, we randomly sampled $N$ samples from the distance matrix 10 times, over

624      an increasing value of *N*. For each set of sampled distances, we computed the minimum observed

625      distance.

626

627      **sOTU novelty.** To assess sOTU novelty, we randomly sampled *N* samples from an sOTU table

628      10 times, over an increasing value of *N*. At each sampling, we computed the number of sOTUs

629      observed with read counts within minimum thresholds. In other words, a minimum threshold of 1

630      is the number of singletons observed in the sampled set, a minimum threshold of 2 is the number

631      of singletons and doubletons, etc.

632

633      **Within-individual beta-diversity.** Many of the individuals in the American Gut Project

634      contributed multiple samples, but at uneven time intervals. In order to explore intrapersonal

635      variation, we replicated the analysis in Lloyd-Price et al. figure 3 (87). Specifically, we

636      determined all time deltas between a subjects samples, and gathered the distributions of beta-

637      diversity between any two samples binned by month. An individual is only represented a single

638      time in a given month, but may be represented in multiple months if they had, for instance,

639      contributed samples over the course of a year.

640

641      **High Performance Liquid Chromatography Mass Spectrometry (HPLC-MS) Analysis.** A

642      total of 498 samples were selected for analysis via mass spectrometry. Specifically, two groups

643      were chosen. First, given the large body of primary literature describing the negative impact of

644      antibiotics on the gut microbiome, and the general interest in this topic from many American Gut

645      participants, we chose 279 samples from individuals (age, BMI, and country matched) who self-

646      reported not having taken antibiotics in the past year, or having taken antibiotics in the past

647     month or week. We chose a second group of 219 samples collected from individuals who

648     answered the question "In an average week, how many different plants do you eat? (e.g., if you

649     consume a can of soup that contains carrots, potatoes and onion, you can count this as 3 different

650     plants; If you consume multi-grain bread, each different grain counts as a plant. Include all fruits

651     in the total)" on the main American Gut Project main survey and who had also completed the

652     VioScreen Food Frequency Questionnaire. When American Gut participants collect samples,

653     they do so on a double headed swab; therefore, all samples chosen for this analysis had one

654     remaining swab head (the first had been used for DNA extraction and microbiome sequencing).

655

656     Cell cultures sample preparation for metabolomics analysis. The supernatant collected from cell

657     cultures (see "expanded bloom assessment" below) were processed to make them compatible

658     with HPLC-MS analysis. The solid phase extraction with wash was carried out to reduce impact

659     of cell culture media, which is highly detrimental for the ESI. The 30 mg sorbent Oasis HLB

660     (Waters, Waltham, MA) SPE cartridges were used to achieve broad metabolite coverage. The

661     cell samples were stored at -80°C and thawed at room temperature immediately prior to

662     extraction. The thawed samples were then centrifuged for 10 minutes at 1200 rpm and extracted.

663     For the SPE extraction, the Oasis HLB SPE cartridge was conditioned with 700μL of 100%

664     HPLC-grade methanol and equilibrated with 700μL of HPLC-grade DI water. The cell

665     supernatant (~350-400μL) was loaded into cartridge and allowed to slowly elute. The loaded

666     SPE wells were then washed with 800μL of 5% methanol in water and the absorbed material was

667     slowly eluted with 200 μL of 100% methanol. Vacuum up to ~ 20 psi was applied for the wells

668     that did not elute within an hour. The collected eluent was stored at -20°C until the HPLC-MS

669     analysis.

28

670

671     <u>Fecal sample preparation for metabolomics analysis.</u> The swab tubes scheduled for analysis were

672     removed from the -80°C freezer and placed on dry ice for the duration of sample processing.

673     Each tube with swab was logged by reading the barcode with barcode scanned and the swab was

674     removed from tube and placed onto a ThermoFisher Scientific (ThermoFisher Scientific,

675     Waltham, MA) 2 ml deep well 96-well plate set on top of dry ice coolant. The top part of each

676     swab's stick was snapped off and discarded. Immediately after filling all of the wells with swabs,

677     200 μL of HPLC-grade 90% v:v ethanol:water solvent was added to each well using

678     multichannel pipette. Four blanks of unused swabs and extraction solvent were included onto

679     each plate. Each plate was then sealed with 96-well plate lid, sonicated for 10 minutes and placed

680     into the refrigerator at 2 °C to extract samples overnight. After extraction, the swabs were

681     removed from wells and discarded, the plates were placed into a lyophilizer, and the entire

682     sample was dried down and then re-suspended in 200 μL 90% v:v ethanol:water. The plates were

683     resealed and centrifuged at 2000 rpm for 10 minutes. The 100 μL aliquots of sample were then

684     transferred onto a Falcon 96-well MS plate using a multichannel pipette, and each plate was

685     immediately sealed with sealing film. The MS plates were centrifuged at 2000 rpm for 10

686     minutes and stored at 2 °C until analysis.

687

688     <u>HPLC-MS analysis.</u> The metabolomics analysis of samples was conducted using reverse phase

689     (RP) high performance liquid chromatography mass spectrometry (HPLC-MS). The HPLC-MS

690     analysis was performed on a Dionex UltiMate 3000 ThermoFisher Scientific high-performance

691     liquid chromatography system (ThermoFisher Scientific, Waltham, MA) coupled to a Bruker

692     impact HD qTOF mass spectrometer. The chromatographic separation was carried out on a

693 Kinetex C18 1.7 μm, 100Å UHPLC column (50 mm x 2.1 mm) (Phenomenex, Torrance, CA),

694 held at 40 ºC during analysis. A total of 5 μL of each sample was injected. Mobile phase A was

695 water, mobile phase B was acetonitrile, both with added 0.1% v:v formic acid. The solvent

696 gradient table was set as follows: initial mobile phase composition was 5% B for 1 min,

697 increased to 40% B over 1 min, then to 100% B over 6 min, held at 100% B for 1 min, decreased

698 back to 5% B in 0.1 min, followed by a washout cycle and equilibration for a total analysis time

699 of 13 min. The scanned m/z range was 80-2000 Th, the capillary voltage was 4500 V, the

700 nebulizer gas pressure was 2 bar, the drying gas flow rate was 9 L/min, and the temperature was

701 200 °C. Each full MS scan was followed by MS/MS using collision-induced dissociation (CID)

702 fragmentation of the seven most abundant ions in the spectrum. For MS/MS, the collision cell

703 collision energy was set at 3 eV and the collision energy was stepped 50%, 75%, 150% and

704 200% to obtain optimal fragmentation for differentially sized ions of different sizes. The scan

705 rate was 3 Hz. A HP-921 lock mass compound was infused during the analysis to carry out post-

706 processing mass correction. All of the raw data are publicly available at the UCSD Center for

707 Computational Mass Spectrometry (*111*) (dataset ID: MassIVE MSV000080179).

708

709 MS data analysis. The collected HPLC-MS raw data files were first converted from Bruker's *d* to

710 mzXML format and then processed with the open source OpenMS 2.0 software (88) in order to

711 deconvolve and align each peak across different chromatograms (feature detection). The

712 alignment window was set at 0.5 minutes, the noise threshold at 1000 counts, the

713 chromatographic peak FWHM value at 20, and the mass error at 30 ppm. All of the peaks that

714 were present in any of the blanks with S/N below 10:1 were removed from the final feature table.

715 The number of features with corresponding MS/MS was as follows: Vioscreen study sample

30

716 cohort: 5144 total MS2 features; antibiotics study samples cohort: 8288 total MS2 features. The

717 number of MS1 features is difficult to estimate exactly as it depends on feature detection settings

718 and the number of samples, but it is typically about 4-5 fold greater depending on the sample.

719 For all of the MS1 features detected across all samples, only ~1-5% are present in an individual

720 sample.

721

722 Chemical annotations were carried out by automatic matching fragmentation spectra to multiple

723 databases using Global Natural Product Social Molecular Networking (GNPS) (89) and then

724 examining the data at the MS/MS level by molecular networking (90). The goal is to retrieve

725 spectra with identical and similar fragmentation patterns and combine them into consensus nodes

726 and clusters, respectively. The consensus node spectra are then compared against public MS/MS

727 libraries to provide molecular annotations (91). Further annotations could be suggested by

728 examining the molecular network (90) (so called propagated annotations). Annotations obtained

729 with precursor and MS/MS matching are considered level two annotations according to the 2007

730 metabolomics standards initiative (92). All molecular networking analysis and annotations are

731 available here: antibiotic use subset (93); types of plants subset (94), cell cultures of isolates (95)

732 and fecal samples co-networked with the cell cultures (96). The raw data contain a significant

733 number of abundant features originating from swab polymers. Therefore, selective background

734 peak removal was carried out specifically for the polymer compounds originating from swabs

735 that were used for the sample collection. The m/z shifts that correspond to the polymer repeating

736 units (44.0262, 88.0524, 132.0786, 176.1049) were identified with GNPS m/z differences

737 frequency plot. The network clusters that contained nodes with the corresponding mass

738 differences were deemed to belong to polymers and all member nodes of the network clusters

31

739    were removed from the feature table (a total of 1632 features/nodes). Principal Coordinates

740    Analysis (PCoA) using a Hellinger distance (97) matrix was used to confirm that the batch effect

741    corresponding to the batches of swabs was mitigated prior to further analysis. To confirm

742    putative annotations, authentic standards were purchased for the linoleic acid (LA; Spectrum

743    Laboratory Products, Inc., USA), conjugated linoleic acid (CLA; mixture of 4 isomers: 9,11 and

744    10,12 isomers, E and Z) (Sigma-Aldrich, USA), and selected antibiotics: tetracycline,

745    oxytetracycline, and doxycyclin (Abcam Inc., USA). For level one identifications, each authentic

746    compound was analyzed under identical experimental conditions and retention time and MS/MS

747    spectra were compared with putatively annotated compounds.

748

749    Selective feature detection. Selective feature extraction was performed with open source

750    MZmine2 software (98). To separate closely eluting LA and CLA isomers as well as separate

751    various N-acyl amides, crop filtering with RT range of 5.4-6.0 minutes and m/z range of 281.246

752    - 281.248 was applied to all chromatograms. Mass detection was performed with a signal

753    threshold of 1.0E2 and a 0.6 s minimum peak width. The mass tolerance was set to 20 ppm and

754    the maximum allowed retention time deviation was set to 5 s. For chromatographic

755    deconvolution, the baseline cutoff algorithm with a 5.0E1 signal threshold was used. The

756    maximum peak width was set to 0.5 min. Similarly, the MS feature for reference compound

757    stercobilin was extracted with a crop filter RT range of 2.0-4.0 minutes and m/z range of

758    595.345-595.355. The stercobilin reference compound was used to assess variability of

759    chromatographic retention times to ensure that the compounds of interest (LA and CLA in

760    particular) retention times were correctly identified. After isotope peak removal, the peak lists of

761    all samples were aligned within the corresponding retention time and mass tolerances. Gap

762    filling was performed on the aligned peak list using the peak finder module with 1% intensity, 10

763    ppm m/z tolerance, and 0.05 min RT tolerance, respectively. After the creation and export of a

764    feature matrix containing the feature retention times, exact mass, and peak areas of the

765    corresponding extracted ion chromatograms, we added sample metadata to the feature matrix

766    metadata of the samples.

767

768    The selective feature extraction with the same settings has been performed for all of the detected

769    compounds listed on the Figure 6A-I ( the m/z range crop filter window was set for

770    corresponding m/z for each compound).

771

772    Molecular Networking. Raw data files were converted to the .mzXML format using Bruker Data

773    Analysis software and uploaded to the GNPS (https://gnps.ucsd.edu/) MassIVE mass

774    spectrometry database (https://massive.ucsd.edu/). Molecular networking was performed to

775    identify spectra shared between different sample types and to identify known molecules in the

776    data set. All annotations are at level 2 according to the proposed minimum standards in

777    metabolomics (92). The data were filtered by removing all MS/MS peaks within +/- 17 Da of the

778    precursor m/z. MS/MS spectra were window-filtered by choosing only the top 6 peaks in the +/-

779    50 Da window throughout the spectrum. The MS spectra were then clustered with MS-Cluster

780    algorithm with a parent mass tolerance of 0.02 Da and a MS/MS fragment ion tolerance of 0.02

781    Da to create consensus spectra (89). Further, consensus spectra that contained less than 4 spectra

782    were discarded. A network was then created where edges were filtered to have a cosine score

783    above 0.65 and more than 5 matched peaks. The edges between two nodes were kept in the

784    network if and only if each of the nodes appeared in each other's respective top 10 most similar

785    nodes. The spectra in the network were then searched against GNPS spectral libraries. The

786    library spectra were filtered in the same manner as the input data. All library matches were

787    required to have a score above 0.7 and at least 6 matched peaks. Molecular networks were

788    visualized and mined using the Cytoscape software (www.cytoscape.org/).

789

790    Molecular networking-based propagation of annotations. The annotation of GPCR agonist

791    compounds was not possible via direct library matching, as their spectra are not present in any

792    MS libraries, but direct comparison with fragmentation patterns presented in (99) allowed us to

793    establish these compounds' identity with level 3 identification (92). Consequently, manual

794    annotation of compounds was carried out in two steps. The exact mass of compounds and their

795    MS/MS fragmentation spectra were matched to the reference spectra found in supplementary

796    info of (99) (fig S4A). Compound m/z 611.5357 was identified in this fashion. In addition,

797    commendamide (330.2640) and its analogue (m/z 344.2799) were identified by matching exact

798    mass of the corresponding ion and by in silico prediction of the MS/MS fragmentation spectra

799    with the CSI:FingerID (100) (fig S4B). For novel molecules that were found within clusters of

800    compounds of interest, but were not described in the literature previously, the structure was

801    postulated using annotation propagation from adjacent annotated nodes in the cluster as

802    described in (89) by assessing differences in parent mass and fragmentation patterns. The key

803    structure, m/z 387.322 has been annotated as N-3-OH-palmitoyl ornithine based on the exact

804    mass and previous annotation (99) as well as analysis of fragmentation pattern to confirm

805    structural moieties of fragments (fig S4C). The rest of the structural assignments have been

806    propagated from that structure. The ornithine moiety has been determined to be present in each

34

807     structure (due to presence of the signature ion with m/z 115.09), and acylation of the hydroxyl is

808     not possible due to insufficient mass of the structures; thus, the changing mass was postulated to

809     correspond to different length of the alkyl substituent (fig 6, in the main text).

810

811     <u>Correlations of Metabolites with Metadata</u>. We have investigated correlations between

812     metabolites (especially those of interest, such as N-Acyl amides) and all of the categories in the

813     metadata. The data were subsetted into the Vioscreen and Antibiotics cohorts and normalized

814     using probabilistic quotient normalization (101).  In order to test the association of the

815     metabolites to the categorical metadata fields we performed the Kruskal–Wallis test followed by

816     Benjamini & Hochberg FDR correction to all metabolites. The significant metabolite-metadata

817     associations ($p$-value adjusted < 0.05) were further connected to GNPS spectral library matches

818     associating the MS1 feature to the MS2 precursor ion in a 10 ppm mass window and 20 seconds

819     retention time window. The results are summarized in table S5.

820

821     <u>Data pretreatment for statistical analysis</u>. A PCoA plot using Hellinger distance (distance matrix:

822     Hellinger; grouping: HCA) was built with all samples in the subset; one sample was found to be

823     an outlier and removed. The data were then filtered to remove features with near-constant, very

824     small values and values with low repeatability using the inter-quartile range estimate. Detailed

825     description of methodology is given in (102). The samples were normalized by sum total of peak

826     intensities, an important step due to large variability of the fecal material load on different swabs.

827     To reduce the effect of background signal and make the sum normalization appropriate, the

828     subtraction of blank and polymer peak features was conducted prior to analysis, as described

829    above. The data were further scaled by mean centering and dividing by standard deviation for

830    each feature.

831

832    The data were split into two groups for downstream analysis. Group one contained samples from

833    individuals answering "More than 30" ($n$=41) and "Less that 10" ($n$=44) to the main American

834    Gut Project survey question "In an average week, how many different plants do you eat?" Group

835    two contained samples from individuals answering "antibiotic use within last week" ($n$=56) and

836    "I have not taken antibiotics in the past year" ($n$=115) to the main American Gut Project survey

837    question "I have taken antibiotics in the last ____." for the Antibiotic history study,

838    correspondingly.

839

840    The resultant features tables were used as input for the Metabonalist software (103). Partial least

841    squares Discriminant Analysis (PLS-DA) (62) was used to explore and visualize variance within

842    data and differences among experimental categories. Random forests (104) (RF) supervised

843    analysis was used to further verify validity of determined discriminating features.

844

845    **Expanded bloom assessment.** The American Gut Project dataset now spans multiple-omics

846    types, and include data that were unavailable during the analysis described in Amir et al. (14). To

847    better understand how the blooming organisms impacted the samples in the American Gut, we 1)

848    performed an additional set of 16S-based experiments; 2) cultured historical samples covering a

849    range of bloom fractions, characterized their metabolites and sequenced the isolates; 3)

850    performed shotgun metagenomics sequencing on the "high bloom" samples; 4) ran the set of

851     samples previously run for HPLC-MS (e.g., the plants and antibiotics cohorts) for shotgun

852     metagenomics, and 5) ran the storage samples from (105) for shotgun metagenomics. The

853     additional sequencing effort was to provide a basis to assess whether functional potential driven

854     by the blooms was impacting any of the biological results discussed in the manuscript. The

855     additional HPLC-MS work was to characterize the metabolites specific to the blooms to remove

856     them from analysis. The additional sequence data generated from the American Gut samples

857     were deposited in EBI under the American Gut accession (ERP012803), and the storage sample

858     data under its accession (ERP015155).

859

860     16S-based bloom experiments. Effect size calculations were computed prior to and following the

861     removal of bloom reads using the procedure described by Amir et al., 2017 (84). The fraction of

862     reads recruiting to blooms was included as a covariate. Effect sizes were assessed over Faith's

863     Phylogenetic Diversity (59), unweighted UniFrac (57) and weighted UniFrac (106). We then

864     computed Pearson and Spearman correlations of the effect sizes, per metric, between the bloom

865     and bloom-removed result (fig 2D, E). In addition to the effect size calculations, we also tested

866     whether the bloom fraction was correlated to any metadata category and did not observe

867     significant correlations.

868

869     We then tested the removal of blooms from other studies in which room temperature shipping

870     was not performed by retrieving a wide variety of human fecal studies from Qiita. UniFrac

871     distance matrices were computed prior to and following bloom removal, followed by Mantel

872     tests. The results of this procedure are outlined in table S4.

873

874      Finally, we correlated the relative intensities of the HPLC-MS data associated with the

875      antibiotics and plants cohorts against the fraction of blooming reads. Critically, we observed a set

876      of spectra that are significantly correlated (table S5) to this fraction. On annotation using

877      molecular networking (discussed in detail the HPLC-MS section), we observed these metabolites

878      to putatively be LysoPE, lysophospholipid (LPL), which has previously been associated with the

879      release of colicin (107). These metabolites were removed from subsequent analyses.

880

881      Culturing. Primary specimens ($n$=214) were selected from three plates based off of the median

882      fraction of reads recruiting to the blooms across the plate, whether the primary specimen still

883      existed, and as to gather samples from at least the US ($n$=116) and UK ($n$=73); additional

884      countries were included in smaller sample sizes and include Australia ($n$=7), Germany ($n$=7),

885      Canada ($n$=3), Croatia ($n$=2), Belgium ($n$=2), France ($n$=1), Austria ($n$=1), Sweden ($n$=1), and

886      the Czech Republic ($n$=1). The bloom typically observed in these samples (and in the full AGP

887      dataset) is an *E. coli* (ID: 04195686f2b70585790ec75320de0d6f from (84)), although a few of

888      the other bloom sequences were represented at high read fraction as well. Samples were retrieved

889      from -80°C and thawed on ice. The swab head was broken off into 500 µl sterile 1x Dulbecco's

890      Phosphate-Buffered Saline and vortexed vigorously for 30 seconds. Serial dilutions from this

891      initial stock were made including 1:10,000 and 1:1,000,000. 10µl of the 1:10,000 dilution were

892      inoculated into 1.5 ml sterile Tryptic Soy broth (TSB, BD cat#2253534) in sterile 96-deep-well

893      plates (community cultures, CC) and incubated overnight at 37°C on an orbital shaker at 500

894      rpm. OD600 values above 0.1 (TSB controls measured ~0.08) were counted as positive growth.

895      Samples with high bloom fraction tended to grow overnight in ambient conditions, samples with

896      a low bloom fraction tended to not grow in these conditions (fig 2A). Additionally, 100 µl of

897 each dilution were plated onto Tryptic Soy agar using sterile glass beads and incubated overnight

898 at 37°C. The following morning, a picture of the best dilution was captured and the most

899 representative colony was selected from each plate and inoculated into 1.5 ml sterile TSB for

900 overnight incubation as above (isolates, IS). The following morning, OD600 measurements were

901 taken and the cultures were pelleted at 3,000 g for 5 min. The supernatant and cell pellets were

902 stored at -20°C for metabolomic analysis and DNA extraction, respectively.

903

904 Shotgun sequencing was performed on all isolates and community cultures using a 1:10

905 miniaturized Nextera library prep with 1 ng gDNA input or up to 1 μl and a 15 cycle PCR

906 amplification. Libraries were quantified with PicoGreen™ dsDNA Assay Kit and 50 ng of each

907 library (or 4 μl maximum) was pooled. The library was size-selected for 200-700 bp using the

908 Sage Bioscience Pippin Prep and sequenced as a paired end 150 cycle run on an Illumina HiSeq

909 2500 v2 in Rapid Run mode at the UCSD IGM Genomics Center. Sequence processing including

910 assembly performed as in the metagenomic processing section below with the exception that "--

911 meta" was not used with SPAdes (108), and read binning against the resulting contigs was not

912 performed. For each isolate, contigs with abnormally high or low coverage as defined by the 1.5

913 × IQR rule were dropped. The characterization of the metabolites from the supernatant using

914 HPLC-MS is discussed in the HPLC-MS section above.

915

916 Following assembly of the draft genomes, taxonomic assessment by Kraken (109) revealed that

917 of the 119 successfully sequenced colony isolate cultures, 95 matched the bloom organisms

918 identified by Amir et al., 2017. Compellingly, 70 of these isolate genomes contained exact 16S

919    sequence matches to a bloom organism identified by (84), including 65 of which matched the

920    dominant *E. coli* bloom in the American Gut (table S4).

921

922    The read data for the isolates were then assessed for predicted biosynthetic gene clusters (BGCs).

923    We used biosyntheticSPAdes (110) to analyze BGCs in the assembly graph of individual

924    genomes. Below we focus on the longest BGCs that are particularly difficult to reconstruct based

925    on ad hoc analysis of contigs and reveal their variations (that likely translate into variations of

926    their natural products). Some of the reconstructed long BGC are ubiquitous (shared by many

927    isolates, albeit with some variations), while others are unique, e.g., present in a single or small

928    number of isolates. We identified BGCs, representing in the alphabet of their domains (table S4),

929    and uncovered variations in their sequence across multiple isolates. Specifically, a ubiquitous

930    BGC similar to the elusive peptide-polyketide genotoxin colibactin and a unique surfactin-like

931    BGC. Colibactin triggers DNA double-strand breaks in eukaryotic cells (111, 112) and induces

932    cellular senescence and metabolic reprogramming in affected mammalian cells (113). Of the 11

933    samples containing the longest colibactin-like BGC, 10 of them contained the exact *E. coli*

934    bloom 16S sequence described above; the 11[th] isolate was actually a canine fecal sample plated

935    alongside human (as the AGP allows participants to submit pet samples).

936

937    Although colibactin is frequently harbored by various *E. coli* strains, the variations of colibactin

938    BGCs across various isolates have not been studied before. Genomic analysis revealed wide

939    variations in colibactin-like BGCs suggesting that various strains produce related but not

940    identical variants of natural products (114). These variations may give rise to the suite of

941    LysoPE-associated spectra identified between the 16S and HPLC-MS datasets.

40

942

943 <u>Shotgun sequencing of the high bloom and storage samples</u>. Previously extracted DNA from the

944 "high bloom" samples used for culturing was obtained, as was previously extracted DNA from

945 Song et al. (105). Shotgun sequencing libraries from a total of 5 ng (or 3.5 μl maximum) gDNA

946 was used in a 1:10 miniaturized KAPA HyperPlus protocol with a 15 cycle PCR amplification.

947 Libraries were quantified with PicoGreen™ dsDNA Assay Kit and 50 ng (or 1 μl maximum) of

948 each library was pooled. The pool was size-selected for 300-700 bp and sequenced as a paired

949 end 150 cycle run on an Illumina HiSeq 2500 v2 in Rapid Run mode at the UCSD IGM

950 Genomics Center. Sequence processing including assembly was performed as in the

951 metagenomic processing section below.

952
953 <u>Functional assessment of conjugated and non-conjugated linoleic acid.</u> To investigate the

954 metabolic potential of gut microbiome for producing conjugated linoleic acid from linoleic acid,

955 we estimated the abundance of linoleic acid isomerase (LAI) in the fecal metagenome. We

956 focused this investigation on the "plants" cohort, which were samples selected to maximize the

957 difference between the number of types of plants metadata category as discussed in the main

958 text. First, we translated the assembled metagenomes to metaproteomes using Prodigal gene

959 prediction software. To map LAI to these metaproteomes, we used a representative LAI protein

960 sequence (UniProt: D2BQ64), which was matched against UniProtKB (via

961 https://www.ebi.ac.uk/Tools/hmmer/) for multiple sequence alignment (MSA). The resulting

962 MSA file in clustal format was then used to generate a hidden Markov model (HMM) profile for

963 LAI using hmmbuild in HMMER software (115). Subsequently, we mapped the resulting HMM

964 profile to sample metaproteomes using hmmsearch with an E-value threshold of 10E-5. We

965 calculated abundances of LAI per sample based on abundance (coverage x length) of LAI

966    containing contigs in each sample, normalized to total sample biomass and performed linear

967    regression between LAI abundances and bloom fraction. We did not note any correlation

968    between metabolic potential of gut metagenome to produce LAI and the fraction of blooming

969    bacteria (samples with no LAI hits were removed from this analysis). Similarly, there was no

970    correlation between CLA abundances and bloom fraction in the samples. These results suggest

971    that our report on the differential abundance of CLA in subjects with different dietary practices

972    (with respect to the number of different types of plants consumed) is unlikely to be confounded

973    by the presence of blooming bacteria.

974

975
976    Storage sample assessment. Metagenomic reads from the storage samples were mapped to the

977    169 isolate assemblies. We then ran model comparison tests on each to determine which

978    mappings were significantly different between frozen samples and samples left out at ambient

979    temperatures for various periods of time.  Using the 'lme' package (116) in R (v3.3.3. R Core

980    Team 2017), linear mixed effects models were applied to the abundances, with individual treated

981    as the random effect.  Mappings were considered to be significantly associated with temperature

982    if the model was significantly improved (ANOVA $p<=0.05$) by incorporating a fixed effect of

983    temperature.  Seven mappings to isolates were found to be significantly increased in samples

984    stored in ambient temperatures compared to frozen samples in both storage studies, of which 3

985    contained the 16S of the dominant *E. coli* bloom in the AGP samples, and 2 contained the 16S

986    from other blooms recognized by (84).

987
988    Shotgun sequence processing. Raw FastQ files were processed using Atropos v1.1.5 (117) to

989    remove adapters and low-quality regions. Putative human genome contaminations were

990    identified and removed by using Bowtie2 v2.3.0 (118) with the "--very-sensitive" option against

991    the human reference genome GRCh37/hg19.

992

993    Sequences were assigned taxonomy using Kraken v1.0.0 (109) against the "standard" database

994    built following the Kraken manual, which contains all complete bacterial, archeal, and viral

995    genomes available from NCBI RefSeq as of Aug. 3, 2017. Results were processed using Bracken

996    v1.0.0 (119) to estimate the relative abundance of species-level taxa.

997

998    Metagenome sequencing data were assembled using SPAdes v3.11.1 (108) with the "−−

999    `meta`" flag enabled. Contigs $\geq$ 1 kb in length were retained and fed to the prokaryotic genome

1000    annotation pipeline Prokka v1.12 (120). putatively individual genomes were inferred using

1001    MaxBin2 v2.2.4 (121).

1002

1003    In parallel, contigs were sheared into 200-bp fragments and taxonomy was assigned using

1004    Kraken (see above). For each contig, the most assigned taxon at each taxonomic rank and the

1005    proportion of sequences assigned to it was inferred.

1006

1007    A total of 3725 genome bins were identified from 677 out of 780 AGP metagenomes, with 5.50

1008    $\pm$ 4.05 bins per sample, and a maximum bin number of 30. Bins with completeness < 50% were

1009    dropped, leaving 1029 bins from 464 samples (2.22 $\pm$ 1.97 bins per sample, maximum bins =

1010    19).

1011

1012    **Filtering Bacterial Blooms for Metabolomics Analysis.** To assess and account for the impact

1013    of the metabolites contributed by these organisms, we have performed HPLC-MS analysis of

1014    cultures of blooming organisms to establish possible contributions, as described above. The raw

43

1015    data are publicly available at the UCSD Center for Computational Mass Spectrometry

1016    (*http://massive.ucsd.edu/*, dataset ID: MassIVE MSV000081777). It was found that there is a

1017    negligible overlap of the bloom-associated metabolites with the compounds detected in AGP

1018    samples (fig 2B). Furthermore, we have verified that none of the compounds discussed in this

1019    work (LA, CLA, compounds on Fig 6A-I) are present in these bloom cultures. The main

1020    organism implicated in bloom was determined to be *E. coli*, as described earlier and MS data

1021    corroborate these findings (fig 2C).

1022

1023    Considering that the metabolites resulting from microbial activity in cultures can differ

1024    significantly from those in vivo (e.g. many of the metabolites could originate not from de novo

1025    synthesis, but rather from microbial modifications of external compounds that are not present in

1026    media, e.g. from the host), we also explored associations of metabolites in AGP metabolomics

1027    samples and blooms. Spearman rank correlation analysis of the fraction of 16S reads in a sample

1028    reporting as bloom to metabolites observed in the same samples revealed several features that

1029    correlate significantly (table S5). There exists a significant overlap between the Antibiotics and

1030    Vioscreen studies subsets, indicating potential common origin of these features.  The strongest

1031    correlation was found for the feature m/z 480.3106 with multiple bloom organisms ($\rho^{\wedge}2 > 0.25$

1032    for *E. coli* at $p < $ 1e-40). This feature was found to also significantly correlate with the principle

1033    coordinates of the PCoA, with and without blooms in the UniFrac matrices for both subsets. The

1034    tentative annotation of this feature is lysoPE, a lysophospholipid (LPL). The LPLs production in

1035    vivo is a result of phospholipase A enzymatic activity associated with Gram-negative bacteria. It

1036    is known that lysoPE is essential for release of colicin (107). Colicin (by itself not detectable

1037    with the MS methodology in this study due to very high molecular mass) is a bacteriocin related

1038    to microbial warfare and is known to be produced by *E. coli*, the major bloomer in AGP. It can

1039    be suggested that the blooming of an organism is related to attempting to kill competitors to

1040    maximize nutrient availability. Importantly, removal of all of the features associated with bloom

1041    does not alter the metabolomics results at all, which indicates that all of the observed biological

1042    trends reported here are not related to blooms.

1043

1044    **Mental health in the American Gut Project.** From AGP cohort, we selected subjects who

1045    endorsed a mental health disorder (depression, schizophrenia, PTSD, and/or bipolar disorder).

1046    This resulted in 1,140 subjects. 636 subjects endorsed at least one of the exclusion criteria

1047    (antibiotic use in the last year, IBD, *C. difficile* infection, pregnancy, Alzheimer's, anorexia or

1048    bulimia, history of substance use disorder, epilepsy or seizure disorder, kidney disease,

1049    phenylketonuria). Out of the remaining 504 subjects, 319 did not provide information regarding

1050    country of residence, hence forming a case cohort of 185 subjects. The remaining samples were

1051    further filtered down to 125 samples to include only high quality fecal microbiome data (at least

1052    1,250 sequences/sample) at a single time point per subject. For those cases, we created a 1:1

1053    matched sample of patients and non-psychiatric comparison (NC) participants based on age ($\pm 5$

1054    years), BMI, history of diabetes, smoking frequency, country of residence, census region (if in

1055    US), and sequencing plate. For each of the cohorts we calculated beta-diversity distance matrices

1056    using Bray-Curtis dissimilarity and weighted UniFrac. On resulting matrices we ran pairwise

1057    PERMANOVA with 999 permutations between "cases" (people who reported mental illness)

1058    and NCs (out matched control dataset). Differential abundance testing was performed using

1059    permutive mean difference test at 10,000 permutations, with discrete FDR (122) correction at

1060    alpha=0.1.

1061

1094

## References

1096  1.  Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A,

1097  Gibbons SM, Ackermann G, Navas-Molina JA, Janssen S, Kopylova E, Vázquez-Baeza Y,

1098  González A, Morton JT, Mirarab S, Zech Xu Z, Jiang L, Haroon MF, Kanbar J, Zhu Q, Jin

1099  Song S, Kosciolek T, Bokulich NA, Lefler J, Brislawn CJ, Humphrey G, Owens SM,

1100  Hampton-Marcell J, Berg-Lyons D, McKenzie V, Fierer N, Fuhrman JA, Clauset A, Stevens

1101  RL, Shade A, Pollard KS, Goodwin KD, Jansson JK, Gilbert JA, Knight R, Earth

1102  Microbiome Project Consortium. 2017. A communal catalogue reveals Earth's multiscale

1103  microbial diversity. Nature 551:457–463.

1104  2.  Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, Kurilshikov A, Bonder MJ,

1105  Valles-Colomer M, Vandeputte D, Tito RY, Chaffron S, Rymenans L, Verspecht C, De

1106  Sutter L, Lima-Mendez G, D'hoe K, Jonckheere K, Homola D, Garcia R, Tigchelaar EF,

1107  Eeckhaudt L, Fu J, Henckaerts L, Zhernakova A, Wijmenga C, Raes J. 2016. Population-

1108        level analysis of gut microbiome variation. Science 352:560–564.

1109   3.   Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez

1110       F, Yamada T, Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H, Xie Y,

1111       Tap J, Lepage P, Bertalan M, Batto J-M, Hansen T, Le Paslier D, Linneberg A, Nielsen HB,

1112       Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu H, Yu C, Li S, Jian M, Zhou Y, Li

1113       Y, Zhang X, Li S, Qin N, Yang H, Wang J, Brunak S, Doré J, Guarner F, Kristiansen K,

1114       Pedersen O, Parkhill J, Weissenbach J, MetaHIT Consortium, Bork P, Ehrlich SD, Wang J.

1115       2010. A human gut microbial gene catalogue established by metagenomic sequencing.

1116       Nature 464:59–65.

1117   4.   Zhernakova A, Kurilshikov A, Bonder MJ, Tigchelaar EF, Schirmer M, Vatanen T, Mujagic

1118       Z, Vila AV, Falony G, Vieira-Silva S, Wang J, Imhann F, Brandsma E, Jankipersadsing SA,

1119       Joossens M, Cenit MC, Deelen P, Swertz MA, LifeLines cohort study, Weersma RK,

1120       Feskens EJM, Netea MG, Gevers D, Jonkers D, Franke L, Aulchenko YS, Huttenhower C,

1121       Raes J, Hofker MH, Xavier RJ, Wijmenga C, Fu J. 2016. Population-based metagenomics

1122       analysis reveals markers for gut microbiome composition and diversity. Science 352:565–

1123       569.

1124   5.   Human Microbiome Project Consortium. 2012. Structure, function and diversity of the

1125       healthy human microbiome. Nature 486:207–214.

1126   6.   Yatsunenko T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M, Magris

1127       M, Hidalgo G, Baldassano RN, Anokhin AP, Heath AC, Warner B, Reeder J, Kuczynski J,

1128       Caporaso JG, Lozupone CA, Lauber C, Clemente JC, Knights D, Knight R, Gordon JI.

1129       2012. Human gut microbiome viewed across age and geography. Nature 486:222–227.

1130   7.   Flores GE, Caporaso JG, Henley JB, Rideout JR, Domogala D, Chase J, Leff JW,

1131    Vázquez-Baeza Y, Gonzalez A, Knight R, Dunn RR, Fierer N. 2014. Temporal variability is

1132    a personalized feature of the human microbiome. Genome Biol 15:531.

1133    8.   Song SJ, Lauber C, Costello EK, Lozupone CA, Humphrey G, Berg-Lyons D, Caporaso JG,

1134    Knights D, Clemente JC, Nakielny S, Gordon JI, Fierer N, Knight R. 2013. Cohabiting family

1135    members share microbiota with one another and with their dogs. Elife 2:e00458.

1136    9.   Smith B, Chu LK, Smith TC, Amoroso PJ, Boyko EJ, Hooper TI, Gackstetter GD, Ryan

1137    MAK, Millennium Cohort Study Team. 2008. Challenges of self-reported medical conditions

1138    and electronic medical records among members of a large military cohort. BMC Med Res

1139    Methodol 8:37.

1140    10.  Morton JT, Sanders J, Quinn RA, McDonald D, Gonzalez A, Vázquez-Baeza Y, Navas-

1141    Molina JA, Song SJ, Metcalf JL, Hyde ER, Lladser M, Dorrestein PC, Knight R. 2017.

1142    Balance Trees Reveal Microbial Niche Differentiation. mSystems 2.

1143    11.  Clemente JC, Pehrsson EC, Blaser MJ, Sandhu K, Gao Z, Wang B, Magris M, Hidalgo G,

1144    Contreras M, Noya-Alarcón Ó, Lander O, McDonald J, Cox M, Walter J, Oh PL, Ruiz JF,

1145    Rodriguez S, Shen N, Song SJ, Metcalf J, Knight R, Dantas G, Dominguez-Bello MG.

1146    2015. The microbiome of uncontacted Amerindians. Sci Adv 1.

1147    12.  Obregon-Tito AJ, Tito RY, Metcalf J, Sankaranarayanan K, Clemente JC, Ursell LK, Zech

1148    Xu Z, Van Treuren W, Knight R, Gaffney PM, Spicer P, Lawson P, Marin-Reyes L, Trujillo-

1149    Villarroel O, Foster M, Guija-Poma E, Troncoso-Corzo L, Warinner C, Ozga AT, Lewis CM.

1150    2015. Subsistence strategies in traditional societies distinguish gut microbiomes. Nat

1151    Commun 6:6505.

1152    13.  Lozupone C, Knight R. 2005. UniFrac: a new phylogenetic method for comparing microbial

1153    communities. Appl Environ Microbiol 71:8228–8235.

1154    14. Amir A, McDonald D, Navas-Molina JA, Debelius J, Morton JT, Hyde E, Robbins-Pianka A,

1155         Knight R. 2017. Correcting for Microbial Blooms in Fecal Samples during Room-

1156         Temperature Shipping. mSystems 2.

1157    15. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP,

1158         Thompson LR, Hyde ER, Gonzalez A, Knight R. 2017. Deblur Rapidly Resolves Single-

1159         Nucleotide Community Sequence Patterns. mSystems 2.

1160    16. Song SJ, Amir A, Metcalf JL, Amato KR, Xu ZZ, Humphrey G, Knight R. 2016. Preservation

1161         Methods Differ in Fecal Microbiome Stability, Affecting Suitability for Field Studies.

1162         mSystems 1.

1163    17. Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, Kurilshikov A, Bonder MJ,

1164         Valles-Colomer M, Vandeputte D, Tito RY, Chaffron S, Rymenans L, Verspecht C, De

1165         Sutter L, Lima-Mendez G, D'hoe K, Jonckheere K, Homola D, Garcia R, Tigchelaar EF,

1166         Eeckhaudt L, Fu J, Henckaerts L, Zhernakova A, Wijmenga C, Raes J. 2016. Population-

1167         level analysis of gut microbiome variation. Science 352:560–564.

1168    18. Zhernakova A, Kurilshikov A, Bonder MJ, Tigchelaar EF, Schirmer M, Vatanen T, Mujagic

1169         Z, Vila AV, Falony G, Vieira-Silva S, Wang J, Imhann F, Brandsma E, Jankipersadsing SA,

1170         Joossens M, Cenit MC, Deelen P, Swertz MA, LifeLines cohort study, Weersma RK,

1171         Feskens EJM, Netea MG, Gevers D, Jonkers D, Franke L, Aulchenko YS, Huttenhower C,

1172         Raes J, Hofker MH, Xavier RJ, Wijmenga C, Fu J. 2016. Population-based metagenomics

1173         analysis reveals markers for gut microbiome composition and diversity. Science 352:565–

1174         569.

1175    19. Ley RE, Lozupone CA, Hamady M, Knight R, Gordon JI. 2008. Worlds within worlds:

1176         evolution of the vertebrate gut microbiota. Nat Rev Microbiol 6:776–788.

1177    20.  McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL,

1178          Knight R, Hugenholtz P. 2012. An improved Greengenes taxonomy with explicit ranks for

1179          ecological and evolutionary analyses of bacteria and archaea. ISME J 6:610–618.

1180    21.  Mirarab S, Nguyen N, Warnow T. 2012. SEPP: SATé-enabled phylogenetic placement. Pac

1181          Symp Biocomput 247–258.

1182    22.  Halfvarson J, Brislawn CJ, Lamendella R, Vázquez-Baeza Y, Walters WA, Bramer LM,

1183          D'Amato M, Bonfiglio F, McDonald D, Gonzalez A, McClure EE, Dunklebarger MF, Knight

1184          R, Jansson JK. 2017. Dynamics of the human gut microbiome in inflammatory bowel

1185          disease. Nat Microbiol 2:17004.

1186    23.  McDonald D, Ackermann G, Khailova L, Baird C, Heyland D, Kozar R, Lemieux M,

1187          Derenski K, King J, Vis-Kampen C, Knight R, Wischmeyer PE. 2016. Extreme Dysbiosis of

1188          the Microbiome in Critical Illness. mSphere 1.

1189    24.  David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, Ling AV,

1190          Devlin AS, Varma Y, Fischbach MA, Biddinger SB, Dutton RJ, Turnbaugh PJ. 2014. Diet

1191          rapidly and reproducibly alters the human gut microbiome. Nature 505:559–563.

1192    25.  Smits SA, Leach J, Sonnenburg ED, Gonzalez CG, Lichtman JS, Reid G, Knight R,

1193          Manjurano A, Changalucha J, Elias JE, Dominguez-Bello MG, Sonnenburg JL. 2017.

1194          Seasonal cycling in the gut microbiome of the Hadza hunter-gatherers of Tanzania.

1195          Science 357:802–806.

1196    26.  McDonald D, Robbins-Pianka A, Mann AE, Vrbanac A, Amir A, Frazier A, Gonzalez A,

1197          Tripathi A, Fahimipour AK, Brennen C, Martino C, Lebrilla C, Lozupone C, M. Lewis C Jr,

1198          Raison C, Zhang C, L. Lauber C, Warinner C, A. Lowry C, Callewaert C, Bloss C,

1199          Huttenhower C, Knights D, Willner D, Galzerani DD, Gonzalez DJ, Mills DA, Chopra D,

1200    Gevers D, Berg-Lyons D, D. Sears D, Wendel D, Wolfe E, Lovelace E, R. Hyde E, Pierce

1201    E, TerAvest E, Montassier E, Bolyen E, Bushman FD, Ackermann G, D. Wu G, Church GM,

1202    Rahnavard G, Saxe G, Gogul G, Humphrey G, D. Holscher H, Ugrina I, German JB,

1203    Gregory Caporaso J, Gilbert J, Wozniak JM, Kerr J, Ravel J, Gaffney J, D. Lewis J, Morton

1204    JT, Suchodolski JS, Jansson JK, Hampton-Marcell JT, Bobe J, Leach J, Raes J, L. Green

1205    J, Metcalf JL, H. Chase J, A. Eisen J, Monk J, Navas-Molina JA, C. Clemente J, Petrosino

1206    J, Ladau J, Shorenstein J, Goodrich J, Gauglitz J, Jacobs J, W. Debelius J, Zengler K, S.

1207    Pollard K, Swanson KS, Lewis K, Mayer K, Bittinger K, Goldasich LD, Dillon L, S. Zaramela

1208    L, R. Thompson L, M. Schriml L, Dominguez-Bello MG, Jankowska MM, Blaser M, Jackson

1209    MA, Pirrung M, Minson M, Kurisu M, Ajami N, Sangwan N, Gottel NR, Chia N, Fierer N,

1210    White O, D. Cani P, Wischmeyer P, Gajer P, Strandwitz P, Hugenholtz P, Dorrestein PC,

1211    Kashyap P, Dutton R, Park RS, Xavier RJ, Knight R, R. Dunn R, Mills RH, Krajmalnik-

1212    Brown R, Ley R, M. Owens S, T. Kelley S, Klemmer S, Matamoros S, Peddada S, Mirarab

1213    S, Janssen S, Moorman S, Holmes S, Schwartz T, Eshoo-Anton TW, Vigers T, Spector T,

1214    Kosciolek T, G. Thackray V, Pandey V, Van Treuren W, Fang X, Chen Y, Vázquez-Baeza

1215    Y, Zech Xu Z, Aksenov AA, Melnik AV, Jarmusch A, Geier J, Pevzner P, Meleshko D,

1216    Behsaz B, Reeve N, Silva R, Zhu Q, Kopylova E, Marotz C, Smarr L, Nguyen T, Mohimani

1217    H, Jiang L, Swafford A, Nguyen D, Song SJ, Sanders K, Benitez RAS, Heale AC, V. Jeste

1218    D, Nguyen TT, Brennan C, Abramson M. 2018. movie_s1.mp4. figshare.

1219    27.  Barberán A, Dunn RR, Reich BJ, Pacifici K, Laber EB, Menninger HL, Morton JM, Henley

1220    JB, Leff JW, Miller SL, Fierer N. 2015. The ecology of microscopic life in household dust.

1221    Proc Biol Sci 282.

1222    28.  Ladau J, Sharpton TJ, Finucane MM, Jospin G, Kembel SW, O'Dwyer J, Koeppel AF,

1223    Green JL, Pollard KS. 2013. Global marine bacterial diversity peaks at high latitudes in

1224    winter. ISME J 7:1669–1677.

1225   29. Gophna U, Konikoff T, Nielsen HB. 2017. Oscillospira and related bacteria - From

1226        metagenomic species to metabolic features. Environ Microbiol 19:835–841.

1227   30. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. 2014. The

1228        carbohydrate-active enzymes database (CAZy) in 2013. Nucleic Acids Res 42:D490–5.

1229   31. El Kaoutari A, Armougom F, Gordon JI, Raoult D, Henrissat B. 2013. The abundance and

1230        variety of carbohydrate-active enzymes in the human gut microbiota. Nat Rev Microbiol

1231        11:497–504.

1232   32. Cummings JH, Macfarlane GT. 1991. The control and consequences of bacterial

1233        fermentation in the human colon. J Appl Bacteriol 70:443–459.

1234   33. Watrous J, Roach P, Alexandrov T, Heath BS, Yang JY, Kersten RD, van der Voort M,

1235        Pogliano K, Gross H, Raaijmakers JM, Moore BS, Laskin J, Bandeira N, Dorrestein PC.

1236        2012. Mass spectral molecular networking of living microbial colonies. Proc Natl Acad Sci U

1237        S A 109:E1743–52.

1238   34. Wang M, Carver JJ, Phelan VV, Sanchez LM, Garg N, Peng Y, Nguyen DD, Watrous J,

1239        Kapono CA, Luzzatto-Knaan T, Porto C, Bouslimani A, Melnik AV, Meehan MJ, Liu W-T,

1240        Crüsemann M, Boudreau PD, Esquenazi E, Sandoval-Calderón M, Kersten RD, Pace LA,

1241        Quinn RA, Duncan KR, Hsu C-C, Floros DJ, Gavilan RG, Kleigrewe K, Northen T, Dutton

1242        RJ, Parrot D, Carlson EE, Aigle B, Michelsen CF, Jelsbak L, Sohlenkamp C, Pevzner P,

1243        Edlund A, McLean J, Piel J, Murphy BT, Gerwick L, Liaw C-C, Yang Y-L, Humpf H-U,

1244        Maansson M, Keyzers RA, Sims AC, Johnson AR, Sidebottom AM, Sedio BE, Klitgaard A,

1245        Larson CB, P CAB, Torres-Mendoza D, Gonzalez DJ, Silva DB, Marques LM, Demarque

1246        DP, Pociute E, O'Neill EC, Briand E, Helfrich EJN, Granatosky EA, Glukhov E, Ryffel F,

1247        Houson H, Mohimani H, Kharbush JJ, Zeng Y, Vorholt JA, Kurita KL, Charusanti P, McPhail

1248      KL, Nielsen KF, Vuong L, Elfeki M, Traxler MF, Engene N, Koyama N, Vining OB, Baric R,

1249      Silva RR, Mascuch SJ, Tomasi S, Jenkins S, Macherla V, Hoffman T, Agarwal V, Williams

1250      PG, Dai J, Neupane R, Gurr J, Rodríguez AMC, Lamsa A, Zhang C, Dorrestein K, Duggan

1251      BM, Almaliti J, Allard P-M, Phapale P, Nothias L-F, Alexandrov T, Litaudon M, Wolfender J-

1252      L, Kyle JE, Metz TO, Peryea T, Nguyen D-T, VanLeer D, Shinn P, Jadhav A, Müller R,

1253      Waters KM, Shi W, Liu X, Zhang L, Knight R, Jensen PR, Palsson BO, Pogliano K,

1254      Linington RG, Gutiérrez M, Lopes NP, Gerwick WH, Moore BS, Dorrestein PC, Bandeira N.

1255      2016. Sharing and community curation of mass spectrometry data with Global Natural

1256      Products Social Molecular Networking. Nat Biotechnol 34:828–837.

1257   35. Sumner LW, Amberg A, Barrett D, Beale MH, Beger R, Daykin CA, Fan TW-M, Fiehn O,

1258      Goodacre R, Griffin JL, Hankemeier T, Hardy N, Harnly J, Higashi R, Kopka J, Lane AN,

1259      Lindon JC, Marriott P, Nicholls AW, Reily MD, Thaden JJ, Viant MR. 2007. Proposed

1260      minimum reporting standards for chemical analysis Chemical Analysis Working Group

1261      (CAWG) Metabolomics Standards Initiative (MSI). Metabolomics 3:211–221.

1262   36. Kishino S, Takeuchi M, Park S-B, Hirata A, Kitamura N, Kunisawa J, Kiyono H, Iwamoto R,

1263      Isobe Y, Arita M, Arai H, Ueda K, Shima J, Takahashi S, Yokozeki K, Shimizu S, Ogawa J.

1264      2013. Polyunsaturated fatty acid saturation by gut lactic acid bacteria affecting host lipid

1265      composition. Proc Natl Acad Sci U S A 110:17808–17813.

1266   37. Coakley M, Ross RP, Nordgren M, Fitzgerald G, Devery R, Stanton C. 2003. Conjugated

1267      linoleic acid biosynthesis by human-derived Bifidobacterium species. J Appl Microbiol

1268      94:138–145.

1269   38. Cohen LJ, Esterhazy D, Kim S-H, Lemetre C, Aguilar RR, Gordon EA, Pickard AJ, Cross

1270      JR, Emiliano AB, Han SM, Chu J, Vila-Farres X, Kaplitt J, Rogoz A, Calle PY, Hunter C,

1271      Bitok JK, Brady SF. 2017. Commensal bacteria make GPCR ligands that mimic human

1272    signalling molecules. Nature 549:48–53.

1273    39. Hauser AS, Attwood MM, Rask-Andersen M, Schiöth HB, Gloriam DE. 2017. Trends in

1274    GPCR drug discovery: new agents, targets and indications. Nat Rev Drug Discov 16:829–

1275    842.

1276    40. Cani PD, Knauf C. 2016. How gut microbes talk to organs: The role of endocrine and

1277    nervous routes. Mol Metab 5:743–752.

1278    41. Clarke G, Stilling RM, Kennedy PJ, Stanton C, Cryan JF, Dinan TG. 2014. Minireview: Gut

1279    microbiota: the neglected endocrine organ. Mol Endocrinol 28:1221–1238.

1280    42. Jiang H, Ling Z, Zhang Y, Mao H, Ma Z, Yin Y, Wang W, Tang W, Tan Z, Shi J, Li L, Ruan

1281    B. 2015. Altered fecal microbiota composition in patients with major depressive disorder.

1282    Brain Behav Immun 48:186–194.

1283    43. Lin P, Ding B, Feng C, Yin S, Zhang T, Qi X, Lv H, Guo X, Dong K, Zhu Y, Li Q. 2017.

1284    Prevotella and Klebsiella proportions in fecal microbial communities are potential

1285    characteristic parameters for patients with major depressive disorder. J Affect Disord

1286    207:300–304.

1287    44. Saraceno B. 2002. The WHO World Health Report 2001 on mental health. Epidemiol

1288    Psichiatr Soc 11:83–87.

1289    45. Zheng P, Zeng B, Zhou C, Liu M, Fang Z, Xu X, Zeng L, Chen J, Fan S, Du X, Zhang X,

1290    Yang D, Yang Y, Meng H, Li W, Melgiri ND, Licinio J, Wei H, Xie P. 2016. Gut microbiome

1291    remodeling induces depressive-like behaviors through a pathway mediated by the host's

1292    metabolism. Mol Psychiatry 21:786–796.

1293    46. Weingarden A, González A, Vázquez-Baeza Y, Weiss S, Humphry G, Berg-Lyons D,

1294    Knights D, Unno T, Bobr A, Kang J, Khoruts A, Knight R, Sadowsky MJ. 2015. Dynamic

1295    changes in short- and long-term bacterial composition following fecal microbiota

1296    transplantation for recurrent Clostridium difficile infection. Microbiome 3:10.

1297  47.  Koenig JE, Spor A, Scalfone N, Fricker AD, Stombaugh J, Knight R, Angenent LT, Ley RE.

1298    2011. Succession of microbial consortia in the developing infant gut microbiome. Proc Natl

1299    Acad Sci U S A 108 Suppl 1:4578–4585.

1300  48.  McDonald D, Robbins-Pianka A, Mann AE, Vrbanac A, Amir A, Frazier A, Gonzalez A,

1301    Tripathi A, Fahimipour AK, Brennen C, Martino C, Lebrilla C, Lozupone C, M. Lewis C Jr,

1302    Raison C, Zhang C, L. Lauber C, Warinner C, A. Lowry C, Callewaert C, Bloss C,

1303    Huttenhower C, Knights D, Willner D, Galzerani DD, Gonzalez DJ, Mills DA, Chopra D,

1304    Gevers D, Berg-Lyons D, D. Sears D, Wendel D, Wolfe E, Lovelace E, R. Hyde E, Pierce

1305    E, TerAvest E, Montassier E, Bolyen E, Bushman FD, Ackermann G, D. Wu G, Church GM,

1306    Rahnavard G, Saxe G, Gogul G, Humphrey G, D. Holscher H, Ugrina I, German JB,

1307    Gregory Caporaso J, Gilbert J, Wozniak JM, Kerr J, Ravel J, Gaffney J, D. Lewis J, Morton

1308    JT, Suchodolski JS, Jansson JK, Hampton-Marcell JT, Bobe J, Leach J, Raes J, L. Green

1309    J, Metcalf JL, H. Chase J, A. Eisen J, Monk J, Navas-Molina JA, C. Clemente J, Petrosino

1310    J, Ladau J, Shorenstein J, Goodrich J, Gauglitz J, Jacobs J, W. Debelius J, Zengler K, S.

1311    Pollard K, Swanson KS, Lewis K, Mayer K, Bittinger K, Goldasich LD, Dillon L, S. Zaramela

1312    L, R. Thompson L, M. Schriml L, Dominguez-Bello MG, Jankowska MM, Blaser M, Jackson

1313    MA, Pirrung M, Minson M, Kurisu M, Ajami N, Sangwan N, Gottel NR, Chia N, Fierer N,

1314    White O, D. Cani P, Wischmeyer P, Gajer P, Strandwitz P, Hugenholtz P, Dorrestein PC,

1315    Kashyap P, Dutton R, Park RS, Xavier RJ, Knight R, R. Dunn R, Mills RH, Krajmalnik-

1316    Brown R, Ley R, M. Owens S, T. Kelley S, Klemmer S, Matamoros S, Peddada S, Mirarab

1317    S, Janssen S, Moorman S, Holmes S, Schwartz T, Eshoo-Anton TW, Vigers T, Spector T,

1318    Kosciolek T, G. Thackray V, Pandey V, Van Treuren W, Fang X, Chen Y, Vázquez-Baeza

56

1319      Y, Zech Xu Z, Aksenov AA, Melnik AV, Jarmusch A, Geier J, Pevzner P, Meleshko D,

1320      Behsaz B, Reeve N, Silva R, Zhu Q, Kopylova E, Marotz C, Smarr L, Nguyen T, Mohimani

1321      H, Jiang L, Swafford A, Nguyen D, Song SJ, Sanders K, Benitez RAS, Heale AC, V. Jeste

1322      D, Nguyen TT, Brennan C, Abramson M. 2018. movie_s2.mp4. figshare.

1323    49.  Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, Owens SM,

1324      Betley J, Fraser L, Bauer M, Gormley N, Gilbert JA, Smith G, Knight R. 2012. Ultra-high-

1325      throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. ISME

1326      J 6:1621–1624.

1327    50.  Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ, Fierer

1328      N, Knight R. 2011. Global patterns of 16S rRNA diversity at a depth of millions of

1329      sequences per sample. Proc Natl Acad Sci U S A 108 Suppl 1:4516–4522.

1330    51.  Apprill A, McNally S, Parsons R, Weber L. 2015. Minor revision to V4 region SSU rRNA

1331      806R gene primer greatly increases detection of SAR11 bacterioplankton. Aquat Microb

1332      Ecol 75:129–137.

1333    52.  Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP,

1334      Thompson LR, Hyde ER, Gonzalez A, Knight R. 2017. Deblur Rapidly Resolves Single-

1335      Nucleotide Community Sequence Patterns. mSystems 2.

1336    53.  McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL,

1337      Knight R, Hugenholtz P. 2012. An improved Greengenes taxonomy with explicit ranks for

1338      ecological and evolutionary analyses of bacteria and archaea. ISME J 6:610–618.

1339    54.  Mirarab S, Nguyen N, Warnow T. 2012. SEPP: SATé-enabled phylogenetic placement. Pac

1340      Symp Biocomput 247–258.

1341    55.  Wang Q, Garrity GM, Tiedje JM, Cole JR. 2007. Naive Bayesian classifier for rapid

1342    assignment of rRNA sequences into the new bacterial taxonomy. Appl Environ Microbiol

1343    73:5261–5267.

1344    56. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N,

1345    Peña AG, Goodrich JK, Gordon JI, Huttley GA, Kelley ST, Knights D, Koenig JE, Ley RE,

1346    Lozupone CA, McDonald D, Muegge BD, Pirrung M, Reeder J, Sevinsky JR, Turnbaugh

1347    PJ, Walters WA, Widmann J, Yatsunenko T, Zaneveld J, Knight R. 2010. QIIME allows

1348    analysis of high-throughput community sequencing data. Nat Methods 7:335–336.

1349    57. Lozupone C, Knight R. 2005. UniFrac: a new phylogenetic method for comparing microbial

1350    communities. Appl Environ Microbiol 71:8228–8235.

1351    58. Gao HX, Regier EE, Close KL. 2016. Prevalence of and trends in diabetes among adults in

1352    the United States, 1988-2012. J Diabetes 8:8–9.

1353    59. Faith DP. 1992. Conservation evaluation and phylogenetic diversity. Biol Conserv 61:1–10.

1354    60. Anderson MJ. 2008. A new method for non-parametric multivariate analysis of variance.

1355    Austral Ecol 26:32–46.

1356    61. Morton JT, Sanders J, Quinn RA, McDonald D, Gonzalez A, Vázquez-Baeza Y, Navas-

1357    Molina JA, Song SJ, Metcalf JL, Hyde ER, Lladser M, Dorrestein PC, Knight R. 2017.

1358    Balance Trees Reveal Microbial Niche Differentiation. mSystems 2.

1359    62. Wold S, Sjöström M, Eriksson L. 2001. PLS-regression: a basic tool of chemometrics.

1360    Chemometrics Intellig Lab Syst 58:109–130.

1361    63. Yatsunenko T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M, Magris

1362    M, Hidalgo G, Baldassano RN, Anokhin AP, Heath AC, Warner B, Reeder J, Kuczynski J,

1363    Caporaso JG, Lozupone CA, Lauber C, Clemente JC, Knights D, Knight R, Gordon JI.

1364    2012. Human gut microbiome viewed across age and geography. Nature 486:222–227.

1365    64. Clemente JC, Pehrsson EC, Blaser MJ, Sandhu K, Gao Z, Wang B, Magris M, Hidalgo G,

1366        Contreras M, Noya-Alarcón Ó, Lander O, McDonald J, Cox M, Walter J, Oh PL, Ruiz JF,

1367        Rodriguez S, Shen N, Song SJ, Metcalf J, Knight R, Dantas G, Dominguez-Bello MG.

1368        2015. The microbiome of uncontacted Amerindians. Sci Adv 1.

1369    65. Obregon-Tito AJ, Tito RY, Metcalf J, Sankaranarayanan K, Clemente JC, Ursell LK, Zech

1370        Xu Z, Van Treuren W, Knight R, Gaffney PM, Spicer P, Lawson P, Marin-Reyes L, Trujillo-

1371        Villarroel O, Foster M, Guija-Poma E, Troncoso-Corzo L, Warinner C, Ozga AT, Lewis CM.

1372        2015. Subsistence strategies in traditional societies distinguish gut microbiomes. Nat

1373        Commun 6:6505.

1374    66. Vázquez-Baeza Y, Pirrung M, Gonzalez A, Knight R. 2013. EMPeror: a tool for visualizing

1375        high-throughput microbial community data. Gigascience 2:16.

1376    67. Bray JR, Curtis JT. 1957. An Ordination of the Upland Forest Communities of Southern

1377        Wisconsin. Ecol Monogr 27:325–349.

1378    68. Mantel N. 1967. The detection of disease clustering and a generalized regression

1379        approach. Cancer Res 27:209–220.

1380    69. Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and

1381        Powerful Approach to Multiple Testing. J R Stat Soc 57:289–300.

1382    70. Moran PAP. 1950. Notes on continuous stochastic phenomena. Biometrika 37:17–23.

1383    71. Maurice CF, Knowles SCL, Ladau J, Pollard KS, Fenton A, Pedersen AB, Turnbaugh PJ.

1384        2015. Marked seasonal variation in the wild mouse gut microbiota. ISME J 9:2423–2434.

1385    72. Cramér H. 1946. Mathematical Methods of Statistics (PMS-9).

1386 73. Welch BL. 1947. The generalisation of student's problems when several different

1387　　population variances are involved. Biometrika 34:28–35.

1388 74. Guo W, Sarkar SK, Peddada SD. 2010. Controlling false discoveries in multidimensional

1389　　directional decisions, with applications to gene expression data on ordered categories.

1390　　Biometrics 66:485–492.

1391 75. Cohen J. 1992. A power primer. Psychol Bull 112:155–159.

1392 76. Zhernakova A, Kurilshikov A, Bonder MJ, Tigchelaar EF, Schirmer M, Vatanen T, Mujagic

1393　　Z, Vila AV, Falony G, Vieira-Silva S, Wang J, Imhann F, Brandsma E, Jankipersadsing SA,

1394　　Joossens M, Cenit MC, Deelen P, Swertz MA, LifeLines cohort study, Weersma RK,

1395　　Feskens EJM, Netea MG, Gevers D, Jonkers D, Franke L, Aulchenko YS, Huttenhower C,

1396　　Raes J, Hofker MH, Xavier RJ, Wijmenga C, Fu J. 2016. Population-based metagenomics

1397　　analysis reveals markers for gut microbiome composition and diversity. Science 352:565–

1398　　569.

1399 77. Anderson MJ. 2001. A new method for non-parametric multivariate analysis of variance:

1400　　NON-PARAMETRIC MANOVA FOR ECOLOGY. Austral Ecol 26:32–46.

1401 78. Koenig JE, Spor A, Scalfone N, Fricker AD, Stombaugh J, Knight R, Angenent LT, Ley RE.

1402　　2011. Succession of microbial consortia in the developing infant gut microbiome. Proc Natl

1403　　Acad Sci U S A 108 Suppl 1:4578–4585.

1404 79. Weingarden A, González A, Vázquez-Baeza Y, Weiss S, Humphry G, Berg-Lyons D,

1405　　Knights D, Unno T, Bobr A, Kang J, Khoruts A, Knight R, Sadowsky MJ. 2015. Dynamic

1406　　changes in short- and long-term bacterial composition following fecal microbiota

1407　　transplantation for recurrent Clostridium difficile infection. Microbiome 3:10.

1408 80. McDonald D, Ackermann G, Khailova L, Baird C, Heyland D, Kozar R, Lemieux M,

1409       Derenski K, King J, Vis-Kampen C, Knight R, Wischmeyer PE. 2016. Extreme Dysbiosis of
1410       the Microbiome in Critical Illness. mSphere 1.

1411  81.  Vázquez-Baeza Y, Gonzalez A, Smarr L, McDonald D, Morton JT, Navas-Molina JA, Knight
1412       R. 2017. Bringing the Dynamic Microbiome to Life with Animations. Cell Host Microbe 21:7–
1413       10.

1414  82.  Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A,
1415       Gibbons SM, Ackermann G, Navas-Molina JA, Janssen S, Kopylova E, Vázquez-Baeza Y,
1416       González A, Morton JT, Mirarab S, Zech Xu Z, Jiang L, Haroon MF, Kanbar J, Zhu Q, Jin
1417       Song S, Kosciolek T, Bokulich NA, Lefler J, Brislawn CJ, Humphrey G, Owens SM,
1418       Hampton-Marcell J, Berg-Lyons D, McKenzie V, Fierer N, Fuhrman JA, Clauset A, Stevens
1419       RL, Shade A, Pollard KS, Goodwin KD, Jansson JK, Gilbert JA, Knight R, Earth
1420       Microbiome Project Consortium. 2017. A communal catalogue reveals Earth's multiscale
1421       microbial diversity. Nature 551:457–463.

1422  83.  Smits SA, Leach J, Sonnenburg ED, Gonzalez CG, Lichtman JS, Reid G, Knight R,
1423       Manjurano A, Changalucha J, Elias JE, Dominguez-Bello MG, Sonnenburg JL. 2017.
1424       Seasonal cycling in the gut microbiome of the Hadza hunter-gatherers of Tanzania.
1425       Science 357:802–806.

1426  84.  Amir A, McDonald D, Navas-Molina JA, Debelius J, Morton JT, Hyde E, Robbins-Pianka A,
1427       Knight R. 2017. Correcting for Microbial Blooms in Fecal Samples during Room-
1428       Temperature Shipping. mSystems 2.

1429  85.  David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, Ling AV,
1430       Devlin AS, Varma Y, Fischbach MA, Biddinger SB, Dutton RJ, Turnbaugh PJ. 2014. Diet
1431       rapidly and reproducibly alters the human gut microbiome. Nature 505:559–563.

1432    86.  Human Microbiome Project Consortium. 2012. Structure, function and diversity of the

1433         healthy human microbiome. Nature 486:207–214.

1434    87.  Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, Brady A, Creasy

1435         HH, McCracken C, Giglio MG, McDonald D, Franzosa EA, Knight R, White O, Huttenhower

1436         C. 2017. Strains, functions and dynamics in the expanded Human Microbiome Project.

1437         Nature 550:61–66.

1438    88.  Röst HL, Sachsenberg T, Aiche S, Bielow C, Weisser H, Aicheler F, Andreotti S, Ehrlich H-

1439         C, Gutenbrunner P, Kenar E, Liang X, Nahnsen S, Nilse L, Pfeuffer J, Rosenberger G,

1440         Rurik M, Schmitt U, Veit J, Walzer M, Wojnar D, Wolski WE, Schilling O, Choudhary JS,

1441         Malmström L, Aebersold R, Reinert K, Kohlbacher O. 2016. OpenMS: a flexible open-

1442         source software platform for mass spectrometry data analysis. Nat Methods 13:741–748.

1443    89.  Wang M, Carver JJ, Phelan VV, Sanchez LM, Garg N, Peng Y, Nguyen DD, Watrous J,

1444         Kapono CA, Luzzatto-Knaan T, Porto C, Bouslimani A, Melnik AV, Meehan MJ, Liu W-T,

1445         Crüsemann M, Boudreau PD, Esquenazi E, Sandoval-Calderón M, Kersten RD, Pace LA,

1446         Quinn RA, Duncan KR, Hsu C-C, Floros DJ, Gavilan RG, Kleigrewe K, Northen T, Dutton

1447         RJ, Parrot D, Carlson EE, Aigle B, Michelsen CF, Jelsbak L, Sohlenkamp C, Pevzner P,

1448         Edlund A, McLean J, Piel J, Murphy BT, Gerwick L, Liaw C-C, Yang Y-L, Humpf H-U,

1449         Maansson M, Keyzers RA, Sims AC, Johnson AR, Sidebottom AM, Sedio BE, Klitgaard A,

1450         Larson CB, P CAB, Torres-Mendoza D, Gonzalez DJ, Silva DB, Marques LM, Demarque

1451         DP, Pociute E, O'Neill EC, Briand E, Helfrich EJN, Granatosky EA, Glukhov E, Ryffel F,

1452         Houson H, Mohimani H, Kharbush JJ, Zeng Y, Vorholt JA, Kurita KL, Charusanti P, McPhail

1453         KL, Nielsen KF, Vuong L, Elfeki M, Traxler MF, Engene N, Koyama N, Vining OB, Baric R,

1454         Silva RR, Mascuch SJ, Tomasi S, Jenkins S, Macherla V, Hoffman T, Agarwal V, Williams

1455         PG, Dai J, Neupane R, Gurr J, Rodríguez AMC, Lamsa A, Zhang C, Dorrestein K, Duggan

1456    BM, Almaliti J, Allard P-M, Phapale P, Nothias L-F, Alexandrov T, Litaudon M, Wolfender J-

1457    L, Kyle JE, Metz TO, Peryea T, Nguyen D-T, VanLeer D, Shinn P, Jadhav A, Müller R,

1458    Waters KM, Shi W, Liu X, Zhang L, Knight R, Jensen PR, Palsson BO, Pogliano K,

1459    Linington RG, Gutiérrez M, Lopes NP, Gerwick WH, Moore BS, Dorrestein PC, Bandeira N.

1460    2016. Sharing and community curation of mass spectrometry data with Global Natural

1461    Products Social Molecular Networking. Nat Biotechnol 34:828–837.

1462    90. Watrous J, Roach P, Alexandrov T, Heath BS, Yang JY, Kersten RD, van der Voort M,

1463    Pogliano K, Gross H, Raaijmakers JM, Moore BS, Laskin J, Bandeira N, Dorrestein PC.

1464    2012. Mass spectral molecular networking of living microbial colonies. Proc Natl Acad Sci U

1465    S A 109:E1743–52.

1466    91. Nguyen DD, Melnik AV, Koyama N, Lu X, Schorn M, Fang J, Aguinaldo K, Lincecum TL Jr,

1467    Ghequire MGK, Carrion VJ, Cheng TL, Duggan BM, Malone JG, Mauchline TH, Sanchez

1468    LM, Kilpatrick AM, Raaijmakers JM, Mot RD, Moore BS, Medema MH, Dorrestein PC.

1469    2016. Indexing the Pseudomonas specialized metabolome enabled the discovery of

1470    poaeamide B and the bananamides. Nat Microbiol 2:16197.

1471    92. Sumner LW, Amberg A, Barrett D, Beale MH, Beger R, Daykin CA, Fan TW-M, Fiehn O,

1472    Goodacre R, Griffin JL, Hankemeier T, Hardy N, Harnly J, Higashi R, Kopka J, Lane AN,

1473    Lindon JC, Marriott P, Nicholls AW, Reily MD, Thaden JJ, Viant MR. 2007. Proposed

1474    minimum reporting standards for chemical analysis Chemical Analysis Working Group

1475    (CAWG) Metabolomics Standards Initiative (MSI). Metabolomics 3:211–221.

1476    93. GNPS antibiotic use subset.

1477    https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=9bd16822c8d448f59a03e6cc8f017f43

1478    94. GNPS plants subset.

1479      https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=d26ae082b1154f73ac050796fcaa6bda

1480  95.  GNPS isolate supernatant.

1481      https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=23f0f5e5c70f4163b445de71d086d186

1482  96.  GNPS fecal samples co-networked.

1483      https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=adcfbba9b4ca448f8b2133559b16d954

1484  97.  1909. Neue Begründung der Theorie quadratischer Formen von unendlichvielen

1485      Veränderlichen. J Reine Angew Math 1909.

1486  98.  Pluskal T, Castillo S, Villar-Briones A, Oresic M. 2010. MZmine 2: modular framework for

1487      processing, visualizing, and analyzing mass spectrometry-based molecular profile data.

1488      BMC Bioinformatics 11:395.

1489  99.  Cohen LJ, Esterhazy D, Kim S-H, Lemetre C, Aguilar RR, Gordon EA, Pickard AJ, Cross

1490      JR, Emiliano AB, Han SM, Chu J, Vila-Farres X, Kaplitt J, Rogoz A, Calle PY, Hunter C,

1491      Bitok JK, Brady SF. 2017. Commensal bacteria make GPCR ligands that mimic human

1492      signalling molecules. Nature 549:48–53.

1493  100.  Dührkop K, Shen H, Meusel M, Rousu J, Böcker S. 2015. Searching molecular structure

1494      databases with tandem mass spectra using CSI:FingerID. Proc Natl Acad Sci U S A

1495      112:12580–12585.

1496  101.  Ejigu BA, Valkenborg D, Baggerman G, Vanaerschot M, Witters E, Dujardin J-C,

1497      Burzykowski T, Berg M. 2013. Evaluation of normalization methods to pave the way

1498      towards large-scale LC-MS-based metabolomics profiling experiments. OMICS 17:473–

1499      485.

1500  102.  Hackstadt AJ, Hess AM. 2009. Filtering for increased power for microarray data

64

1501    analysis. BMC Bioinformatics 10:11.

1502    103.    Xia J, Wishart DS. 2016. Using MetaboAnalyst 3.0 for Comprehensive Metabolomics

1503    Data Analysis. Curr Protoc Bioinformatics 55:14.10.1–14.10.91.

1504    104.    Breiman L. 2001. 10.1023/A:1010933404324. Machine Learning.

1505    105.    Song SJ, Amir A, Metcalf JL, Amato KR, Xu ZZ, Humphrey G, Knight R. 2016.

1506    Preservation Methods Differ in Fecal Microbiome Stability, Affecting Suitability for Field

1507    Studies. mSystems 1.

1508    106.    Lozupone CA, Hamady M, Kelley ST, Knight R. 2007. Quantitative and qualitative beta

1509    diversity measures lead to different insights into factors that structure microbial

1510    communities. Appl Environ Microbiol 73:1576–1585.

1511    107.    Pugsley AP, Goldzahl N, Barker RM. 1985. Colicin E2 production and release by

1512    Escherichia coli K12 and other Enterobacteriaceae. J Gen Microbiol 131:2673–2686.

1513    108.    Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,

1514    Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,

1515    Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its

1516    applications to single-cell sequencing. J Comput Biol 19:455–477.

1517    109.    Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification

1518    using exact alignments. Genome Biol 15:R46.

1519    110.    Meleshko D, Mohimani H, Hajirasouliha I, Medema MH, Korobeynikov A, Pevzner P.

1520    BiosyntheticSPAdes: Reconstructing Biosynthetic Gene Clusters From Assembly Graphs.

1521    submitted.

1522    111.    Nougayrède J-P, Homburg S, Taieb F, Boury M, Brzuszkiewicz E, Gottschalk G,

1523    Buchrieser C, Hacker J, Dobrindt U, Oswald E. 2006. Escherichia coli induces DNA double-
1524    strand breaks in eukaryotic cells. Science 313:848–851.

1525    112.    Putze J, Hennequin C, Nougayrède J-P, Zhang W, Homburg S, Karch H, Bringer M-A,
1526    Fayolle C, Carniel E, Rabsch W, Oelschlaeger TA, Oswald E, Forestier C, Hacker J,
1527    Dobrindt U. 2009. Genetic structure and distribution of the colibactin genomic island among
1528    members of the family Enterobacteriaceae. Infect Immun 77:4696–4703.

1529    113.    Secher T, Samba-Louaka A, Oswald E, Nougayrède J-P. 2013. Escherichia coli
1530    producing colibactin triggers premature and transmissible senescence in mammalian cells.
1531    PLoS One 8:e77157.

1532    114.    Gurevich A, Mikheenko A, Shlemov A, Korobeynikov A, Mohimani H, Pevzner P. 2018.
1533    Modification-tolerant database search reveals surprising diversity of peptidic natural
1534    products. Nature Microbiology.

1535    115.    Finn RD, Clements J, Eddy SR. 2011. HMMER web server: interactive sequence
1536    similarity searching. Nucleic Acids Res 39:W29–37.

1537    116.    Bates D, Mächler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models
1538    Using lme4. J Stat Softw 67.

1539    117.    Didion JP, Martin M, Collins FS. 2017. Atropos: specific, sensitive, and speedy trimming
1540    of sequencing reads. PeerJ 5:e3720.

1541    118.    Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat
1542    Methods 9:357–359.

1543    119.    Lu J, Breitwieser FP, Thielen P, Salzberg SL. 2017. Bracken: estimating species
1544    abundance in metagenomics data.

1545 120.   Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics

1546       30:2068–2069.

1547 121.   Wu Y-W, Simmons BA, Singer SW. 2016. MaxBin 2.0: an automated binning algorithm

1548       to recover genomes from multiple metagenomic datasets. Bioinformatics 32:605–607.

1549 122.   Jiang L, Amir A, Morton JT, Heller R, Arias-Castro E, Knight R. 2017. Discrete False-

1550       Discovery Rate Improves Identification of Differentially Abundant Microbes. mSystems 2.

1551 123.   Cohen LJ, Kang H-S, Chu J, Huang Y-H, Gordon EA, Reddy BVB, Ternei MA, Craig JW,

1552       Brady SF. 2015. Functional metagenomic discovery of bacterial effectors in the human

1553       microbiome and isolation of commendamide, a GPCR G2A/132 agonist. Proc Natl Acad Sci

1554       U S A 112:E4825–34.

1555

1556 **Figure 1.** Population characteristics. **(A)** Participants across the world have sent in samples to

1557 the American Gut, although the primary geographic regions of participation are in North

1558 America and the United Kingdom; the report a participant receives is depicted. **(B)** The primary

1559 sample breakdown for subsequent analyses. Red denotes reasons samples were removed. **(C)**

1560 Between the two largest populations, the US ($n$=6,634) and the UK ($n$=2,071), we observe a

1561 significant difference in alpha diversity. **(D)** In a meta-analysis, the largely industrialized

1562 population that makes up the American Gut exhibits significant differential abundances to non-

1563 industrialized populations.

1564

1565 **Figure 2.** Blooms and effect sizes. **(A)** The fraction of 16S reads that recruit to bloom reads

1566 defined by Amir et al. 2017 is strongly associated with the likelihood for microbial growth under

1567 aerobic culture conditions on rich media. **(B)** Molecular network of the metabolites observed in

1568   the supernatant from cultures (*n*=217) derived from fecal samples. The nodes in red (n=239) are

1569   metabolites associated with *E. coli*. **(C)** Overlap of metabolites between AGP samples and

1570   blooms.  **(D)** Unweighted UniFrac effect sizes. The inset shows the correlation of effect sizes

1571   when including or excluding the bloom 16S reads (Pearson *r*=0.91, *p*=3.76x10$^{-57}$). **(E)** Weighted

1572   UniFrac effect sizes. The inset shows the correlation of the effect sizes when including or

1573   excluding bloom 16S reads (Pearson *r*=0.42, *p*=1.71x10$^{-6}$); the outlier is the 16S bloom fraction

1574   of the sample.

1575

1576   **Figure 3.** OTU and beta-diversity novelty. **(A)** The AGP data placed into the context of extant

1577   microbial diversity at a global scale. **(B)** A phylogenetic tree showing the diversity spanned by

1578   the AGP, and the HMP in the context of Greengenes and the EMP. **(C)** sOTU novelty over

1579   increasing numbers of samples in the AGP; the AGP appears to have begun to reach saturation

1580   and is contrasted with **(D)** Yatsunenko et al. 2012 which unlike the AGP had extremely deep

1581   sequencing per sample. **(E)** The minimum observed UniFrac distance between samples over

1582   increasing numbers of samples for the AGP and the HMP; inset is from 0-500 samples. **(F)** An

1583   AGP "trading card" of an sOTU of interest (shown in full in fig S2).

1584

1585   **Figure 4.** Temporal and spatial patterns. **(A)** 565 individuals had multiple samples. Distances

1586   between samples within an individual shown at 1 month, 2 months, etc out to over 1 year;

1587   between subject distances shown in "BSD." Even at one year, the median distance between a

1588   participant's samples is less than the median between participant distance. **(B)** Within the US,

1589   spatial processes of sOTUs appear driven by stochastic processes as few sOTUs exhibit spatial

1590   autocorrelation (Moran's *I*) on the full dataset or partitions (e.g., participants older than 20). **(C)**

1591    Distance-decay relationship for Bray-Curtis dissimilarities between subject pairs that are within

1592    100km (great-circle distance) radius of one another (Mantel test; $r$=0.036, adjusted $p$=0.03). Inset

1593    shows the largest radius (i.e., the contiguous US). Darker colors indicate higher-frequency bins.

1594    Dashed lines represent fits from linear models to raw data. **(D)** Mantel correlogram of estimated

1595    $r$ coefficients, significance of distance-decay relationships, and radius (x-axis). Red points

1596    represent neighborhood sizes that were significant (adjusted $p$-values < 0.05). **(E)** Characterizing

1597    a large bowel resection using the AGP, the EMP, a hunter-gatherer population, and ICU patients

1598    in an unweighted UniFrac principal coordinates plot. A state change was observed in the

1599    resulting microbial community. The change in the microbial community immediately following

1600    surgery is the same as the distance between a marine sediment sample and a plant rhizosphere

1601    sample.

1602

1603    **Figure 5.** Diversity of plants in a diet. **(A)** Procrustes analysis of fecal samples from ($n$=1,596)

1604    individuals using Principal Components of the Vioscreen FFQ responses and Principal

1605    Coordinates of the unweighted UniFrac distances ($M^2$=0.988) colored by diet; Procrustes tests

1606    the fit of one ordination space to another. PCA shows grouping by diets such as Vegan

1607    suggesting self-reported diet type is consistent with differences in micro and macro nutrients as

1608    recorded by the FFQ, however these dietary differences do not explain relationships between the

1609    samples in 16S space. **(B)** The full AGP dataset including skin and oral samples through

1610    unweighted UniFrac and Principal Coordinates Analysis highlighting a lack of apparent

1611    clustering by diet type. **(C)** Dietary conjugated linoleic acid levels as reported by the FFQ

1612    between the extremes of plant diversity consumption, and **(D)** the observed levels of CLA by

1613    HPLC-MS. **(E)** Differential abundances of sOTUs (showing the most specific taxon name per

1614 sOTU) between those who eat fewer than 10 plants per week vs. those who eat over 30 per week.

1615 **(F)** The molecules, linoleic acid (LA) and conjugated linoleic acid (CLA) (only trans-, trans-

1616 isomers are shown) were found to comprise the octadecadienoic acid found to be the key feature

1617 in this difference in number of plants consumption.

1618

1619 **Figure 6.** Molecular novelty in the gut microbiome. **(A-I)** Molecular sub-network of N-acyl

1620 amides. Cluster/nodes of microbially-derived G protein-coupled receptor agonistic molecules

1621 detected in human fecal samples are shown. Molecules B, G and H have been described

1622 (compounds 1, 2 & 4b (38) and commendamide (123)); molecules A, C, D, E and I are

1623 previously not reported (proposed structures are shown). **(J)** Compound occurrence frequency

1624 plot. Examples of compounds originating from food (piperine, black pepper alkaloid), host

1625 (stercobilin, heme catabolism product), bacterial activity (lithocholic acid, microbially-modified

1626 bile acid) or exogenous compounds such as antibiotics (rifaximin) or other drugs (lisinopril, high

1627 blood pressure medication) are shown. **(K-N)** Alpha and beta-diversity assessments of antibiotic

1628 and plants cohorts; insets depict minimum observed beta-diversity over increasing samples.

1629

1630 **Supplementary Text:**

1631  **Effect size comparisons**

1632  **Multi-cohort replication detail**

1633  **Projects using the American Gut infrastructure**

1634  **American Gut Survey**

1635  **Supplemental references**

1636

1637 **Supplementary Figures:**

1638 **Figure S1.** Workflow and population scale analyses. **(A)** Heatmap of income levels from the US

1639 Census and American Gut participant locations. **(B)** Sample flowchart for what sample sets

1640 correspond to each analysis. **(C)** Using PLS-DA we observed separation between US ($n$=6,634)

1641 and UK ($n$=2,071) fecal samples. **(D)** We performed a Principal Coordinates analysis comparing

1642 children over the age of 3 and adults from industrialized ($n$=4,643 AGP samples, $n$=4,927

1643 samples total), remote farming ($n$=131), and hunter-gatherer ($n$=30) lifestyles.

1644

1645 **Figure S2.** Trading cards and LS's samples compared to ICU patients and AGP participants and

1646 diet state change analysis. **(A)** Unweighted UniFrac distance distributions for the sample

1647 immediate prior to surgery vs. all ICU fecal samples, and distances of the sample immediately

1648 following surgery vs. all ICU fecal samples (Kruskal Wallis $H$=79.774, $p$=4.198x$^{-19}$). **(B)** Same

1649 as panel **(A)** except comparing against all AGP fecal samples (Kruskal Wallis $H$=8117.734,

1650 $p$=0.0). **(C)** The median distances of each sample in Larry's longitudinal dataset compared to

1651 both ICU and AGP. The last pre-surgery sample is on day 25 and the first post-surgery sample is

1652 day 27. **(D)** A principal coordinates analysis of UniFrac distances of the American Gut Project,

1653 samples from the "extreme" diet study by David et al. (85), and the Earth Microbiome Project.

1654 No obvious state change by the diet of the participants in David et al. is observed.

1655

1656 **Figure S3.** Dietary levels of linoleic acid based on validated food frequency questionnaire

1657 responses, and the detected linoleic acid by mass spectrometry did not differ significantly

1658 between groups consuming few or many types of plants per week.

1659

1660     **Figure S4**. Metabolomic identification and annotation. **(A)** Manual annotation via comparison of

1661     experimental MS fragmentation patterns to those given in (99). Top panel: reference spectrum

1662     for the "Compound 2" in (99); bottom panel: experimental MS/MS spectrum for the parent ion

1663     m/z 611.5357. The compound is annotated as 3-(myristoyloxy)palmitoyl lysine. **(B)** I*n silico*

1664     annotation using CSI:FingerID (100) for the ion with m/z 330.2640. Top panel: experimental

1665     fragmentation pattern explained by the putative fragmentation tree; bottom panel: the possible

1666     candidate structures ranked by match %. The top structure with 71.02% match corresponds to

1667     commendamide. **(C)** Manual annotation via comparison of experimental exact mass to that of

1668     identified compound in (100), N-3-OH-palmitoyl ornithine. The peaks in experimental MS/MS

1669     spectrum are examined and compared to theoretical fragments that would result from breaking

1670     bonds in the proposed structure. The structure is deemed to be consistent with the N-3-OH-

1671     palmitoyl ornithine annotation.

1672

1673     **Supplemental Tables:**

1674     **Table S1.** Summary of sample numbers and type in the American Gut other studies, sample

1675     distributions by country and territory, sample distributions by US state, US participant

1676     demographics and per sequencing round sample accessions in EBI.

1677

1678     **Table S2.** American Gut data dictionary, proportion of responses per AG survey question that

1679     are represented as a single question; multiselect responses were omitted as these are stored in the

1680     metadata as per response type, informal dietary questions and correlations to the food frequency

1681     questionnaire, effect size results without bloom sOTUs, variable mapping with Falony et al. 2016

1682     Science.

1683

1684    **Table S3.** sOTUs relevant to the balance analyses, and summary of differentially abundant taxa

1685    in UK cohort (negative effect size indicated the taxon is more prevalent in control (NC)

1686    subjects).

1687

1688    **Table S4.** Application of the filter for blooms to other human fecal studies which were not

1689    subjected to room temperature shipping, taxonomy of the draft isolate genomes, the specific

1690    bloom 16S sOTUs observed, and ubiquitous colibactin-like biosynthetic gene clusters (top) and a

1691    unique surfactin-like biosynthetic gene cluster observed in the bloom isolates.

1692

1693    **Table S5.** A set of molecular features which appeared to significantly correlate to the bloom

1694    fraction, and Kruskal–Wallis tests for metabolites in the Antibiotics and Vioscreen cohorts of

1695    samples.

1696

1697

A

B

C

D

A

PC2 (8.81%)

| | |
|---|---|
| ■ | Omnivore |
| ■ | Vegetarian |
| ■ | Omnivore but do not eat red meat |
| ■ | Vegetarian but eat seafood |
| ■ | Vegan |

PC1 (19.9%)

B

PC2 (5.21%)

PC1 (15.04%)

C

Conjugated Linoleic Acid (g)

****

Less than 10 plants    More than 30 plants

D

Conjugated Linoleic Acid Abundance (a.u.)

Less than 10    More than 30

E

Less than 10 plants    More than 30 plants

Lachnospiraceae
Lachnospiraceae
Ruminococcaceae
Ruminococcaceae
Ruminococcaceae
*Blautia*
*Oscillospira*
Clostridiales
Clostridiales
*F. prausnitzii*
Erysipelotrichaceae
Lachnospiraceae
Lachnospiraceae
*Blautia*
*Blautia*
*Dorea*
*Ruminococcus*
*R. gnavus*
*E. lenta*
*R. torques*

0.02    0.00    0.02

Mean (proportion)

F

Linoleic Acid

Conjugated Linoleic Acid

A — NH$_2$  m/z 415.353  Proposed structure  H$^+$  HO  O  OH  $C_{14}H_{29}$

D — NH$_2$  m/z 413.337  Proposed structure  H$^+$  HO  O  OH  $C_{14}H_{27}$

B — NH$_2$  m/z 387.322  Compound 4b  H$^+$  HO  O  OH  $C_{12}H_{25}$

E — NH$_2$  m/z 359.290  Proposed structure  H$^+$  HO  O  OH  $C_{10}H_{21}$

C — NH$_2$  m/z 401.338  Proposed structure  H$^+$  HO  O  OH  $C_{13}H_{27}$

F — NH$_2$  m/z 373.306  Proposed structure  H$^+$  HO  O  OH  $C_{11}H_{23}$

Network nodes:
415.354
387.322
413.337
359.290
401.338
373.306

Edge labels: −28.032, −26.0, −54.047, −28.032, −14.016, −14.016, −14.016, −28.03

G — NH$_2$  m/z 611.536  Compound 2  H$^+$  HO  O  O  O  0.0  611.536

H — m/z 330.264  Compound 1  Commendamide  H$^+$  HO  O  OH  0.0  330.264

I — H$^+$  O  OH  m/z 344.280  Proposed structure  0.0  344.280  $C_{14}H_{29}$

J —
Number of samples
350
300
250
200
150
100
50

Specific node number
1  2  4  8  16  32  64  128  256  512  1024  2048  4096  8192

Stercobilin
Piperine
Lithocholic acid
Lisinopril
Rifaximin

K —
Observed molecular features
180
160
140
120
100
80
60
40
Antibiotic use
● No use in the past year
● Within the last month
Bootstrap depth
100  500  900  1300  1700  2100  2500  2900  3300  3700  4100  4500  4900  5300  5700  6100  6500  6900  7300  7700  8100  8500  8900  9300  9700
(inset) Minimum Bray Curtis Distance — Sample count

L —
Observed 16S sOTUs
180
160
140
120
100
80
60
40
Bootstrap depth
100  500  900  1300  1700  2100  2500  3300  3700
(inset) Minimum Bray Curtis Distance — Sample count

M —
Observed molecular features
200
175
150
125
100
75
50
Plants consumed per week
● More than 30
● Less than 10
Bootstrap depth
100  500  900  1300  1700  2100  2500  2900  3300  3700  4100  4500  4900  5300  5700  6100  6500  6900  7300  7700  8100  8500  8900  9300  9700
(inset) Minimum Bray Curtis Distance — Sample count

N —
Observed 16S sOTUs
200
180
160
140
120
100
80
60
40
Bootstrap depth
100  500  900  1300  1700  2100  2500  3300  3700
(inset) Minimum Bray Curtis Distance — Sample count