

1

2 **A comprehensive human minimal gut metagenome extends the host's**
3 **metabolic potential**

4

5 Marcos Parras-Moltó¹, Daniel Aguirre de Cárcer^{1*}

6 ¹Departamento de Biología, Universidad Autónoma de Madrid, Madrid, 28049 Spain.
7 (mparmol@gmail.com, daniel.aguirre@uam.es).

8

9

10

11

12

13

14

15

16

17

18 ***Correspondence to:** Daniel Aguirre de Cárcer García. Darwin 2, 28049 Madrid,
19 Spain. Phone: +34914978429. Fax: +34914978344 daniel.aguirre@uam.es

20

21

22 **ABSTRACT**

23 Accumulating evidence suggests that humans should be considered as holobionts in
24 which the gut microbiota plays essential functions. Initial metagenomic studies reported
25 a pattern of shared genes in the gut microbiome of different individuals, leading to the
26 definition of the minimal gut metagenome as the set of microbial genes necessary for
27 homeostasis, and present in all healthy individuals. Despite its interest, this concept has
28 received little attention following its initial description in terms of various ubiquitous
29 pathways in Western cohorts. This study analyzes the minimal gut metagenome of the
30 most comprehensive dataset available, including individuals from agriculturalist and
31 industrialist societies, also embodying highly diverse ethnic and geographical
32 backgrounds. The outcome, based on metagenomic predictions for community
33 composition data, resulted in a minimal metagenome comprising 3,412 gene clusters,
34 mapping to 1,856 reactions and 128 metabolic pathways predicted to occur across all
35 individuals. These results were substantiated by the analysis of two additional datasets
36 describing the microbial community compositions of larger Western cohorts, as well as
37 a substantial shotgun metagenomics dataset. Subsequent analyses showed the plausible
38 metabolic complementarity provided by the minimal gut metagenome to the human
39 genome.

40

41 **Keywords:** Human gut Microbiome, 16S rRNA gene, PICRUSTs, Community
42 Assembly, Metagenomics.

43

44

45 The study of the human gut microbiome has drawn from different disciplines (e.g.
46 microbiology, ecology, genomics), and has substantiated the idea that humans should be
47 considered as holobionts (1) in which the gut microbiota plays essential functions (2, 3).
48 Knowledge of what constitutes a healthy gut microbiome is regarded as pivotal (4) for
49 the development of predictive models for diagnosis and management of gut
50 microbiome-related maladies. However, the strong inter-subject variability in
51 community composition observed in cross-sectional studies (5) hindered an early
52 definition of a set of bacterial species common to all healthy humans (6). While, recent
53 efforts have been able to detect such a health-related set in terms of shared taxonomic
54 assignments (4, 7), and more precisely in terms of shared 16S sequence clusters of
55 varying phylogenetic depth (8), the idea that a healthy gut microbiome ‘core’ may exist
56 only in terms of function (9) remains widespread.

57 In this regard, early high-throughput shotgun metagenomic studies already reported a
58 strong pattern of shared genes in the gut microbiome of different individuals (10, 11).
59 These results led to the definition of a novel concept; the minimal gut metagenome (11),
60 defined as the set of microbial genes necessary for the homeostasis of the whole gut
61 ecosystem, and expected to be present in all healthy humans. The idea that the gut
62 microbiome provides a specific set of functionalities shared by all individuals is
63 intuitive. However, it is still unclear whether these functionalities could arise from a
64 shared set of genes or from different combinations of genes. Moreover, if the host were
65 to play a greatly diminished role as a selective force on its resident gut microbiome,
66 when compared to external factors such as diet, then there would be no set of microbial
67 functionalities shared by all humans.

68 Nevertheless, despite its potential as a conceptual framework with which to study the
69 gut ecosystem, the minimal gut metagenome concept has received little attention in the

70 literature following its initial definition and description in terms of various ubiquitous
71 metabolic pathways (9-11) and recent description of prevalent pathways in a larger
72 cohort (12).

73 Hence, the aim of the present study is to recapitulate the minimal human gut
74 metagenome conceptual framework, and provide a proof-of-concept of its utility. More
75 specifically, we set out to identify the ‘core genes’ (defined as the set of genes detected
76 in all individuals) , jointly comprising the minimal gut metagenome, as well as the ‘core
77 reactions’ (defined as the set of metabolic reactions detected in all individuals).
78 According to the minimal gut metagenome concept, the former should be related to gut
79 homeostasis at large (i.e. not only metabolic homeostasis). On the other hand,
80 knowledge on the latter should improve our understanding of the gut microbiome's
81 ability to augment human metabolism.

82

83 For knowledge of the minimal gut metagenome to be most useful, it should pertain
84 more to *Homo sapiens* as a species, and hence should not be solely focused on Western
85 cohorts. Unfortunately, most human gut shotgun metagenomic datasets are very
86 restricted in terms of lifestyles and ethnicities, mostly arising from Western and(or)
87 industrialist cohorts (9-13).

88 In this study, 16S rRNA gene-based metagenomic predictions were employed in the
89 assessment of the minimal human gut metagenome to be able to profit from the more
90 comprehensive 16S datasets. These datasets greatly outclass available human gut
91 shotgun metagenomic datasets in terms of cohort size, geographic distribution, ethnic
92 and lifestyle diversity, and to a certain extent depth of sequencing. In a sense, one read
93 in a shotgun metagenomics dataset represents one gene count, while one read in a 16S

94 amplicon survey represents, *via* metagenomic prediction, one genome count. However,
95 the use of metagenomic predictions presents various limitations and possible biases,
96 which have been explored previously (14), the most noteworthy being that it only infers
97 the bacterial and archaeal component of the metagenome, is significantly affected by
98 both the quality of available genome annotations and the fact that available genomes are
99 not evenly distributed across the phylogeny, or the lack of perfect one-to-one mapping
100 between genomes and even full-length 16S sequences. Nevertheless, the ability to count
101 almost three orders of magnitude more genes in a metagenomic sample per sequence
102 (with the number of bacterial genes per genome normally in the very few thousand),
103 even as a prediction, is still useful. In this study, functional predictions based on 16S
104 phylogenetic marker gene sequences were obtained using PICRUSt, a computational
105 approach which has shown large and significant correlation in predicting metagenomic
106 abundances from 16S measurements (Spearman $r = 0.82$, $p < 0.001$) and synthetic
107 communities (Spearman $r = 0.9$, $p < 0.001$)(14). To date, PICRUSt has been used in a
108 myriad of scientific works and different research scenarios, such as the analysis of
109 environmental samples (15), medically-relevant communities (16), or *in vitro*
110 assemblies (17). This study analyzes the minimal gut metagenome of the most
111 comprehensive dataset available (dataset *Global*: 382 individuals from rural Malawi,
112 metropolitan U.S.A., and Venezuelan Amerindians(18). See **Table 1**)), which, despite
113 its comparatively smaller cohort size, is far more inclusive in terms of global
114 distribution, lifestyle, and ethnicity, specifically including agriculturalist, and
115 industrialist societies from three continents.

116 We compare the Global dataset with two larger Western cohorts (dataset *Flemish*: 873
117 individuals from Belgium (4); and dataset *Twins*: 2,727 individuals from U.K. (19)), as
118 well as to a substantial shotgun metagenomics dataset (Dataset *Shotgun*: KEGG

119 Orthology identifiers (KOs) (20) abundances from 123 individuals from U.S.A.,
120 Europe, and China. Obtained from Bradley and Pollard 2017 (21)), and compared with
121 the human genome to assess the degree to which the minimal metagenome may
122 complement and expand its host's metabolic potential.

123

124 **RESULTS**

125 The authors of the PICRUSt paper state that there is a significant negative correlation
126 (Spearman $r = -0.4$, $P < 0.001$) between NSTI values and Spearman correlation
127 between empirical shotgun metagenome abundances and PICRUSt predictions based on
128 16S sequences.(14). Here, NSTI values for the different sample sets of *Global*
129 (0.135 ± 0.021 , 0.098 ± 0.018 , and 0.131 ± 0.023 for Malawian, U.S.A., and Venezuelan
130 samples, respectively; see **Suppl. Fig 1**) were lower (generally correlated with higher
131 correlation between metagenomic measurements and 16S predictions) than those
132 previously reported for soil samples (0.17 ± 0.02) which showed a significant [$P < 0.001$]
133 correlation between predictions and matched shotgun metagenomics assignments (14).
134 Also, the more extreme NSTI values reported for the Human Microbiome Project
135 dataset, with NSTI values ranging 0.10-0.15, still presented high correlation coefficients
136 between metagenomic measurements and 16S predictions (14).

137 The results show that 5,865 KO groups were predicted as present in *Global's* pan-
138 metagenome, while the minimal gut metagenome represented 3,412 KOs (i.e. core
139 genes), which can in turn be mapped to 1,856 reactions (i.e. core reactions) and 128
140 complete metabolic pathways (**Additional file 1**).

141 As could be expected, lowering the prevalence threshold used to define core reactions
142 (100%) increased the number of core reactions, but mainly in a gentle-slope linear

143 fashion (**Suppl. Fig. 2**). The core metagenome was very similar among the three
144 distinct sample sets comprising *Global* (**Figure 1A**), with U.S.A.'s set showing the
145 smallest set of core reactions, and less overlap with Malawian and Venezuelan samples.
146 On the other hand, *Global*'s core reaction set was comparatively similar to those
147 obtained using Western-like datasets *Twins* and *Flemish* (**Figure 1B**).

148 The presented core reactions were predicted from 16S profiles using an ancestral-state
149 reconstruction algorithm (PICRUSt). However, the set of core reactions was
150 substantiated by the use of Tax4fun (22), a taxonomy assignments-based approach
151 (**Figure 1C**). PICRUSt's predictions seem conservative (more appropriate for a
152 minimum estimate, as intended) since they are a subset of Tax4fun predictions. More
153 importantly, *Global*'s core reaction set presented a high overlap to that obtained from

154 a substantial shotgun metagenomics dataset targeting the human gut microbiome (21),
155 chosen among those publicly available based on the number of individuals and
156 geographic and ethnic distribution (**Figure 1D**). The 463 reactions described as core in
157 *Shotgun* but not in *Global* (**Figure 1D**) likely arise from the smaller size of the
158 *Shotgun*'s cohort as well as its increased lifestyle, environmental and genetic
159 homogeneity (**Table 1**). On the other hand, the great majority of core reactions in
160 *Global* not described as core in *Shotgun* still presented a very high prevalence in the
161 dataset (**Suppl. Fig. 3**); 1,735 out of 1,856 (93.5%) core reactions in *Global* are also
162 core reactions (100% prevalence) in *Shotgun*. Only 37 (2%) core reactions in *Global*
163 have a prevalence level < 95% in *Shotgun*, and 6 (0.32%) reactions have a prevalence
164 level below 75%. No apparent shared functional or taxonomic origin affiliation was
165 found for these six reactions. Within the *Global* dataset, there was a positive correlation
166 between prevalence and average abundance (**Suppl. Fig. 4**). Nevertheless, while all core

167 reactions featured relatively high average abundance values, many similarly abundant
168 reactions presented lower prevalence values.

169 In addition to providing an improved description of the human minimal gut
170 metagenome, the present study aimed at assessing its complementarity to the human
171 genome. In this regard, the metabolic complementarity judged by the Metabolic
172 Complementarity Index (23) was >2 times larger when considering the human
173 metabolism being complemented by *Global's* minimal gut metagenome, when
174 compared to the inverse (0.0807 and 0.0386, respectively).

175 Considering two metabolites as linked if they represent the substrate and product of a
176 core reaction, within the overall metabolic map (**Figure 2, Suppl. Fig. 5**) 199 microbial
177 metabolites link with 89 *Homo sapiens* metabolites through 256 core reactions,
178 representing the predicted extended metabolic capability of the human holobiont
179 provided by its gut ecosystem. Additionally, the map pinpoints 55 core reactions and 84
180 metabolites with no apparent connection to *Homo sapiens* metabolism, as well as 36
181 core reactions able to link *Homo sapiens* metabolites by reactions different to those
182 carried-out by enzymes encoded within the human genome.

183 Not surprisingly, several core reactions are implicated in the production of short-chain
184 fatty acids (SCFAs), such as butyrate and acetate, which are known to have an active
185 role in normal human physiology (e.g. fuel for several cell types, regulation of gene
186 expression, differentiation, and inflammation) (24, 25). Another hallmark of the
187 predicted minimal gut metagenome relates to the presence of core reactions implicated
188 in the production of several vitamins (B1, B2, B5, B6, B9, H, K1, K2, L1, coenzyme
189 B12), several of which had previously been shown to be produced by common gut
190 commensals (26).

191

192

193 **DISCUSSION**

194 The NSTI values that we obtained for human gut microbiome samples fall within the
195 range of NSTI values for samples in the PICRUSt validation that had high correlation
196 between metagenomic abundance measurements and 16S predictions.(14). In this
197 regard, an enhanced and updated report on the utility, correlation between predicted and
198 experimental measurements, and accuracy of PICRUSt's predictions would be
199 welcomed by the community, more so since this area of development seems to remain
200 active (27, 28). The values obtained were not homogenous among the three distinct
201 sample sets in *Global*, with values for both the Venezuelan and Malawian samples
202 being roughly 35% higher than that of the U.S.A. samples. In this regard, the detected
203 functional overlap could somewhat be inflated since the reference genome set employed
204 is likely biased towards strains obtained from industrialist countries.

205 Interestingly, the results indicate that the U.S.A. population restricted the number of
206 detected core reactions, since Venezuela and Malawian samples presented an additional
207 156 reactions with 100% prevalence in their joint dataset, compared to <20 exclusively
208 shared with 100% prevalence between U.S.A. samples and any of the other groups.
209 Moreover, these values may be conservative, since the reference genomes may be
210 biased towards bacterial strains more frequent in industrialist countries. This reduction
211 in functional overlap provides circumstantial support to the emerging concern that
212 industrialist populations may have lost the microbial diversity needed to adequately
213 sustain a healthy host (29).

214 The results presented herein are influenced by the fact that the metagenomic prediction
215 approach employed is, to a certain extent, biased, as explained before. As such, the core
216 genes and reactions reported should be taken cautiously. Thus, validation of each
217 particular core reaction in the ecosystem, as well as the possibility of each core
218 metabolite traversing the membrane, along with its potential significance to the host, is
219 beyond the scope of this study. Nevertheless, returning to the three possible scenarios of
220 shared functionality in the human gut pan-microbiome postulated above; i) no shared
221 functionality, ii) shared functionality related to different combinations of genes, and iii)
222 shared functionality related to a shared combination of genes, the results are strongly
223 supportive of the latter. Thus, we believe that the minimal gut metagenome idea indeed
224 represents a potentially useful conceptual framework with which to improve our
225 knowledge of the role played by the human gut microbiome on maintaining host
226 homeostasis.

227 The results also indicate that the human gut minimal metagenome may extensively
228 contribute to the human holobiont's metabolic potential. The core reactions reported
229 here represent a highly restrictive set, since reactions need to be present in all subjects to
230 achieve the 'core' status. Most importantly, these core reactions were predicted as
231 present in all subjects from a cohort including individuals from agriculturalist and
232 industrialist societies, also embodying highly diverse genetic, ethnic, and geographical
233 backgrounds. Furthermore, the results were validated using additional large-cohort
234 datasets, as well as a substantial shotgun metagenomics dataset. Hence, the described
235 minimal gut metagenome now pertains more to *Homo sapiens* as a species, rather than
236 to industrialist societies of particular ethnic and geographical backgrounds. Finally, our
237 results seem to indicate that the minimal metagenome has a greater role in
238 complementing the human metabolism than the other way around.

239

240 **MATERIALS AND METHODS**

241 **Datasets.** All datasets comprised 16S rRNA gene sequences obtained using primer pair
242 F515-R806 targeting the V4 hypervariable region, with the exception of dataset *Shotgun*
243 which included KOs abundances obtained through shotgun sequencing of metagenomic
244 DNA (21). All sequence data was derived from stool samples from healthy subjects
245 over 3 years old, with no history of recent antibiotic treatment prior to sampling (see
246 **Table 1**).

247 **Metagenomic predictions.** QIIME (30) scripts were employed during initial sequence
248 processing (Additional file 2). Briefly, datasets were independently processed as
249 follows; first subsampled to the minimum common depth. Then, chimeric sequences
250 were identified with *usearch61* (31) and removed. Finally, sequences were clustered
251 into OTUs using Greengenes (32) 0.97 representative sequence dataset (May 2013) as
252 reference using *usearch61*. Subsequently, PICRUSt scripts were employed to first
253 normalize OTU abundances by 16S rRNA gene copy number, and then transform
254 normalized OTU abundances into KO abundances. Correlation between predictions and
255 measurements was evaluated using NSTI as a proxy for the Spearman coefficient, as
256 they are strongly negatively and significantly correlated (14). Tax4Fun (22), an
257 alternative metagenome prediction pipeline, was also employed with *Global* dataset
258 following the suggested standard procedure.

259 Since more than one KO group may carry out a particular reaction, KO abundances
260 were mapped to KEGG reactions. In cases where a KO mapped to more than one
261 reaction, all reactions linked to the KO were scored. KOs and reactions appearing in all

262 individuals in the datasets were defined as ‘core’. Finally, the MinPath algorithm (33)
263 was used for biological pathway reconstruction from core KOs.

264 **Metabolic complementarity assessments.** Host-microbiome cooperation was assessed
265 with NetCooperate (23) using the Metabolic Complementarity Index. This index
266 provides a quantification of the extent to which two species may support one another
267 through biosynthetic complementarity. There is no threshold for ‘complementarity’ and
268 ‘no complementarity’, and hence the metrics have to be employed in a comparative
269 manner (23). Here, the index was used to study both moieties of the human holobiont;
270 the human genome and the minimal gut metagenome. Hence, the reciprocal analysis
271 evaluates the relative strength of each moiety complementing the other. To do so, core
272 reactions were transformed into linked KEGG compounds, and then analyzed with
273 NetCooperate. To further assess such complementarity, both the core reactions and the
274 reactions encoded by the human genome were imported into the interactive metabolic
275 pathway explorer iPATH3.0 (34).

276

277 **Availability of data and material**

278 The datasets analyzed during the current study are available from their original source
279 (as stated above). Core KOs, Reactions and Compounds are available within **Additional**
280 **file 1**. Additional intermediate result files and scripts are available from the
281 corresponding author on request for research purposes.

282 **Competing interests**

283 The authors declare that they have no competing interests.

284 **Funding**

285 This work was funded by the Spanish Ministry of Science and Innovation grant
286 BIO2016-80101-R.

287 **Authors' contributions**

288 DA Conceived the idea and wrote the manuscript. MP and DA analyzed the datasets.

289 **Acknowledgements**

290 We thank the Bioinformatics Unit at CBMSO for their support.

291

292 **REFERENCES**

- 293 1. **Bordenstein SR, Theis KR.** 2015. Host Biology in Light of the Microbiome: Ten
294 Principles of Holobionts and Hologenomes. *PLOS Biology* **13**:e1002226.
- 295 2. **Backhed F, Ley RE, Sonnenburg JL, Peterson DA, Gordon JI.** 2005. Host-bacterial
296 mutualism in the human intestine. *Science* **307**:1915-1920.
- 297 3. **Gill SR, Pop M, Deboy RT, Eckburg PB, Turnbaugh PJ, Samuel BS, Gordon JI, Relman
298 DA, Fraser-Liggett CM, Nelson KE.** 2006. Metagenomic analysis of the human distal
299 gut microbiome. *Science* **312**:1355-1359.
- 300 4. **Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, Kurilshikov A, Bonder
301 MJ, Valles-Colomer M, Vandeputte D, Tito RY, Chaffron S, Rymenans L, Verspecht C,
302 De Sutter L, Lima-Mendez G, D'hoel K, Jonckheere K, Homola D, Garcia R, Tigchelaar
303 EF, Eeckhaut L, Fu J, Henckaerts L, Zhernakova A, Wijmenga C, Raes J.** 2016.
304 Population-level analysis of gut microbiome variation. *Science* **352**:560-564.
- 305 5. **Aguirre de Cárcer D, Cuiv PO, Wang T, Kang S, Worthley D, Whitehall V, Gordon I,
306 McSweeney C, Leggett B, Morrison M.** 2011. Numerical ecology validates a
307 biogeographical distribution and gender-based effect on mucosa-associated bacteria
308 along the human colon. *ISME J* **5**:801-809.
- 309 6. **Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI.** 2007. The
310 human microbiome project. *Nature* **449**:804-810.
- 311 7. **Zhang J, Guo Z, Xue Z, Sun Z, Zhang M, Wang L, Wang G, Wang F, Xu J, Cao H, Xu H, Lv
312 Q, Zhong Z, Chen Y, Qimuge S, Menghe B, Zheng Y, Zhao L, Chen W, Zhang H.** 2015. A
313 phylo-functional core of gut microbiota in healthy young Chinese cohorts across
314 lifestyles, geography and ethnicities. *ISME J* **9**:1979-1990.
- 315 8. **Aguirre de Cárcer D.** 2018. The human gut pan-microbiome presents a compositional
316 core formed by discrete phylogenetic units. *Scientific Reports* **8**:14069.
- 317 9. **Turnbaugh PJ, Hamady M, Yatsunenkov T, Cantarel BL, Duncan A, Ley RE, Sogin ML,
318 Jones WJ, Roe BA, Affourtit JP, Egholm M, Henrissat B, Heath AC, Knight R, Gordon JI.**
319 2009. A core gut microbiome in obese and lean twins. *Nature* **457**:480-484.
- 320 10. **Human Microbiome Project Consortium.** 2012. Structure, function and diversity of the
321 healthy human microbiome. *Nature* **486**:207-214.
- 322 11. **Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N,
323 Levenez F, Yamada T, Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H,
324 Xie Y, Tap J, Lepage P, Bertalan M, Batto JM, Hansen T, Le Paslier D, Linneberg A,**

- 325 **Nielsen HB, Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu H, Yu C, Jian M,**
326 **Zhou Y, Li Y, Zhang X, Qin N, Yang H, Wang J, Brunak S, Dore J, Guarner F, Kristiansen**
327 **K, Pedersen O, Parkhill J, Weissenbach J, Bork P, Ehrlich SD.** 2010. A human gut
328 microbial gene catalogue established by metagenomic sequencing. *Nature* **464**:59-65.
- 329 12. **Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, Brady A, Creasy**
330 **HH, McCracken C, Giglio MG, McDonald D, Franzosa EA, Knight R, White O,**
331 **Huttenhower C.** 2017. Strains, functions and dynamics in the expanded Human
332 Microbiome Project. *Nature* **550**:61.
- 333 13. **Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, Arumugam M, Kultima JR, Prifti E,**
334 **Nielsen T, Juncker AS, Manichanh C, Chen B, Zhang W, Levenez F, Wang J, Xu X, Xiao**
335 **L, Liang S, Zhang D, Zhang Z, Chen W, Zhao H, Al-Aama JY, Edris S, Yang H, Hansen T,**
336 **Nielsen HB, Brunak S, Kristiansen K, Guarner F, Pedersen O, Dore J, Ehrlich SD, Bork**
337 **P.** 2014. An integrated catalog of reference genes in the human gut microbiome. *Nat*
338 *Biotechnol* **32**:834-841.
- 339 14. **Langille MG, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,**
340 **Burkepile DE, Vega Thurber RL, Knight R, Beiko RG, Huttenhower C.** 2013. Predictive
341 functional profiling of microbial communities using 16S rRNA marker gene sequences.
342 *Nat Biotechnol* **31**:814-821.
- 343 15. **Bier RL, Voss KA, Bernhardt ES.** 2015. Bacterial community responses to a gradient of
344 alkaline mountaintop mine drainage in Central Appalachian streams. *ISME J* **9**:1378-
345 1390.
- 346 16. **Buffie CG, Bucci V, Stein RR, McKenney PT, Ling L, Gobourne A, No D, Liu H,**
347 **Kinnebrew M, Viale A, Littmann E, van den Brink MR, Jenq RR, Taur Y, Sander C,**
348 **Cross JR, Toussaint NC, Xavier JB, Pamer EG.** 2015. Precision microbiome
349 reconstitution restores bile acid mediated resistance to *Clostridium difficile*. *Nature*
350 **517**:205-208.
- 351 17. **Goldford JE, Lu N, Bajić D, Estrela S, Tikhonov M, Sanchez-Gorostiaga A, Segrè D,**
352 **Mehta P, Sanchez A.** 2018. Emergent simplicity in microbial community assembly.
353 *Science* **361**:469-474.
- 354 18. **Yatsunenkov T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M,**
355 **Magris M, Hidalgo G, Baldassano RN, Anokhin AP, Heath AC, Warner B, Reeder J,**
356 **Kuczynski J, Caporaso JG, Lozupone CA, Lauber C, Clemente JC, Knights D, Knight R,**
357 **Gordon JI.** 2012. Human gut microbiome viewed across age and geography. *Nature*
358 **486**:222-227.
- 359 19. **Goodrich JK, Davenport ER, Beaumont M, Jackson MA, Knight R, Ober C, Spector TD,**
360 **Bell JT, Clark AG, Ley RE.** 2016. Genetic Determinants of the Gut Microbiome in UK
361 Twins. *Cell Host Microbe* **19**:731-743.
- 362 20. **Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M.** 2016. KEGG as a
363 reference resource for gene and protein annotation. *Nucleic Acids Res* **44**:17.
- 364 21. **Bradley PH, Pollard KS.** 2017. Proteobacteria explain significant functional variability in
365 the human gut microbiome. *Microbiome* **5**:017-0244.
- 366 22. **Asshauer KP, Wemheuer B, Daniel R, Meinicke P.** 2015. Tax4Fun: predicting
367 functional profiles from metagenomic 16S rRNA data. *Bioinformatics* **31**:2882-2884.
- 368 23. **Levy R, Carr R, Kreimer A, Freilich S, Borenstein E.** 2015. NetCooperate: a network-
369 based tool for inferring host-microbe and microbe-microbe cooperation. *BMC*
370 *Bioinformatics* **16**:015-0588.
- 371 24. **Donohoe DR, Garge N, Zhang X, Sun W, O'Connell TM, Bunker MK, Bultman SJ.** 2011.
372 The microbiome and butyrate regulate energy metabolism and autophagy in the
373 mammalian colon. *Cell Metab* **13**:517-526.
- 374 25. **Furusawa Y, Obata Y, Fukuda S, Endo TA, Nakato G, Takahashi D, Nakanishi Y, Uetake**
375 **C, Kato K, Kato T, Takahashi M, Fukuda NN, Murakami S, Miyauchi E, Hino S, Atarashi**
376 **K, Onawa S, Fujimura Y, Lockett T, Clarke JM, Topping DL, Tomita M, Hori S, Ohara O,**

- 377 **Morita T, Koseki H, Kikuchi J, Honda K, Hase K, Ohno H.** 2013. Commensal microbe-
378 derived butyrate induces the differentiation of colonic regulatory T cells. *Nature*
379 **504**:446-450.
- 380 26. **Letunic I, Yamada T, Kanehisa M, Bork P.** 2008. iPath: interactive exploration of
381 biochemical pathways and networks. *Trends Biochem Sci* **33**:101-103.
- 382 27. **Douglas GM, Beiko RG, Langille MGI.** 2018. Predicting the Functional Potential of the
383 Microbiome from Marker Genes Using PICRUSt. *Methods Mol Biol*:8728-8723_8711.
- 384 28. **Douglas GM, Maffei VJ, Zaneveld J, Yurgel SN, Brown JR, Taylor CM, Huttenhower C,**
385 **Langille MGI.** 2019. PICRUSt2: An improved and extensible approach for metagenome
386 inference. *bioRxiv*:672295.
- 387 29. **Blaser MJ.** 2018. The Past and Future Biology of the Human Microbiome in an Age of
388 Extinctions. *Cell* **172**:1173-1177.
- 389 30. **Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer**
390 **N, Pena AG, Goodrich JK, Gordon JL, Huttley GA, Kelley ST, Knights D, Koenig JE, Ley**
391 **RE, Lozupone CA, McDonald D, Muegge BD, Pirrung M, Reeder J, Sevinsky JR,**
392 **Turnbaugh PJ, Walters WA, Widmann J, Yatsunencko T, Zaneveld J, Knight R.** 2010.
393 QIIME allows analysis of high-throughput community sequencing data. *Nat Meth*
394 **7**:335-336.
- 395 31. **Edgar RC.** 2010. Search and clustering orders of magnitude faster than BLAST.
396 *Bioinformatics* **26**:2460-2461.
- 397 32. **DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D,**
398 **Hu P, Andersen GL.** 2006. Greengenes, a chimera-checked 16S rRNA gene database
399 and workbench compatible with ARB. *Appl Environ Microbiol* **72**:5069-5072.
- 400 33. **Ye Y, Doak TG.** 2009. A parsimony approach to biological pathway
401 reconstruction/inference for genomes and metagenomes. *PLoS Comput Biol* **5**:14.
- 402 34. **Darzi Y, Letunic I, Bork P, Yamada T.** 2018. iPath3.0: interactive pathways explorer v3.
403 *Nucleic Acids Res* **46**:W510-W513.

404

405 **Additional file 1.**

406 Excel file (.xlsx)

407 Core KOs, Reactions, Compounds and Pathways.

408

409 **Additional file 2.**

410 Word file (.docx)

411 QIIME scripts employed.

412

413

414

415

416

417

418 **FIGURE LEGENDS**

419

420 **Figure 1. Venn diagrams depicting the overlap in core reactions between different**
421 **datasets and software. Panel A:** Different sample sets within *Global*. Values refer to
422 the analysis with the same number of individuals per population (50). **Panel B:**
423 Different 16S datasets. Values refer to the analysis with the same number of sequences
424 per sample (8,000). **Panel C:** Differences between metagenomic prediction software.
425 **Panel D:** Differences between 16S (*Global*) and shotgun metagenomics (*Shotgun*)
426 datasets.

427 **Figure 2. The minimal gut metagenome extends human metabolic potential.** Nodes
428 in the map correspond to chemical compounds and edges represent enzymatic reactions.
429 The figure provides an iPath2.0 representation of KEGG metabolic pathways, where
430 reactions catalyzed by enzymes encoded in the human genome appear in blue, while
431 core reactions of the human gut pan-microbiome not encoded also by the human
432 genome, appear in red.

433

434 **Supplementary Figure 1. Distribution of NSTI values among the three sample sets**
435 **in *Global*.**

436 **Supplementary Figure 2. The number of core reactions varies with prevalence**
437 **threshold.** [Linear regression; $y = -8.57 + 2899x$, $R^2 = 0.92$]

438 **Supplementary Figure 3. Prevalence of *Global* core reactions in *Shotgun*.** Dots
439 represent all reactions detected in *Shotgun*. Their prevalence in the dataset is recorded
440 along the y-axis, and those reactions with 100% prevalence in *Global* (core) appear in a
441 different color.

442 **Supplementary Figure 4. Prevalence Vs. average abundance values in *Global*.** Dots
443 represent all reactions predicted in the dataset, core reactions depicted in red.

444 **Supplementary Figure 5. The gut metagenome extends human metabolic potential.**
445 Nodes in the map correspond to chemical compounds and edges represent enzymatic
446 reactions. The figure provides an iPath2.0 representation of KEGG metabolic pathways,
447 where reactions catalyzed by enzymes encoded in the human genome appear in blue,
448 while reactions of the human gut pan-microbiome not encoded also by the human
449 genome appear in either red (100% prevalence), orange (50% prevalence), or yellow
450 (1% prevalence).

451

452

453

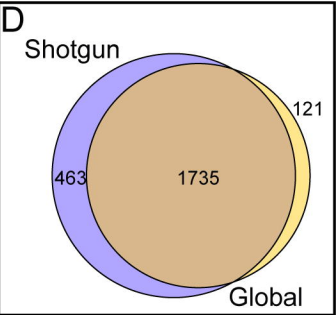
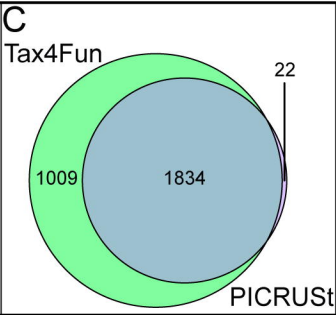
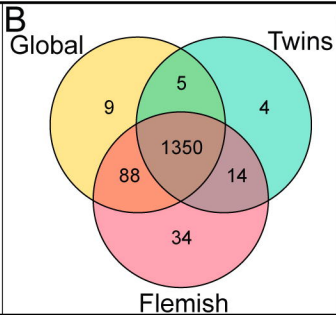
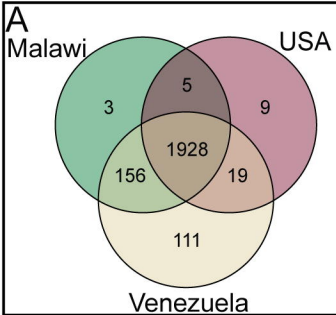
454

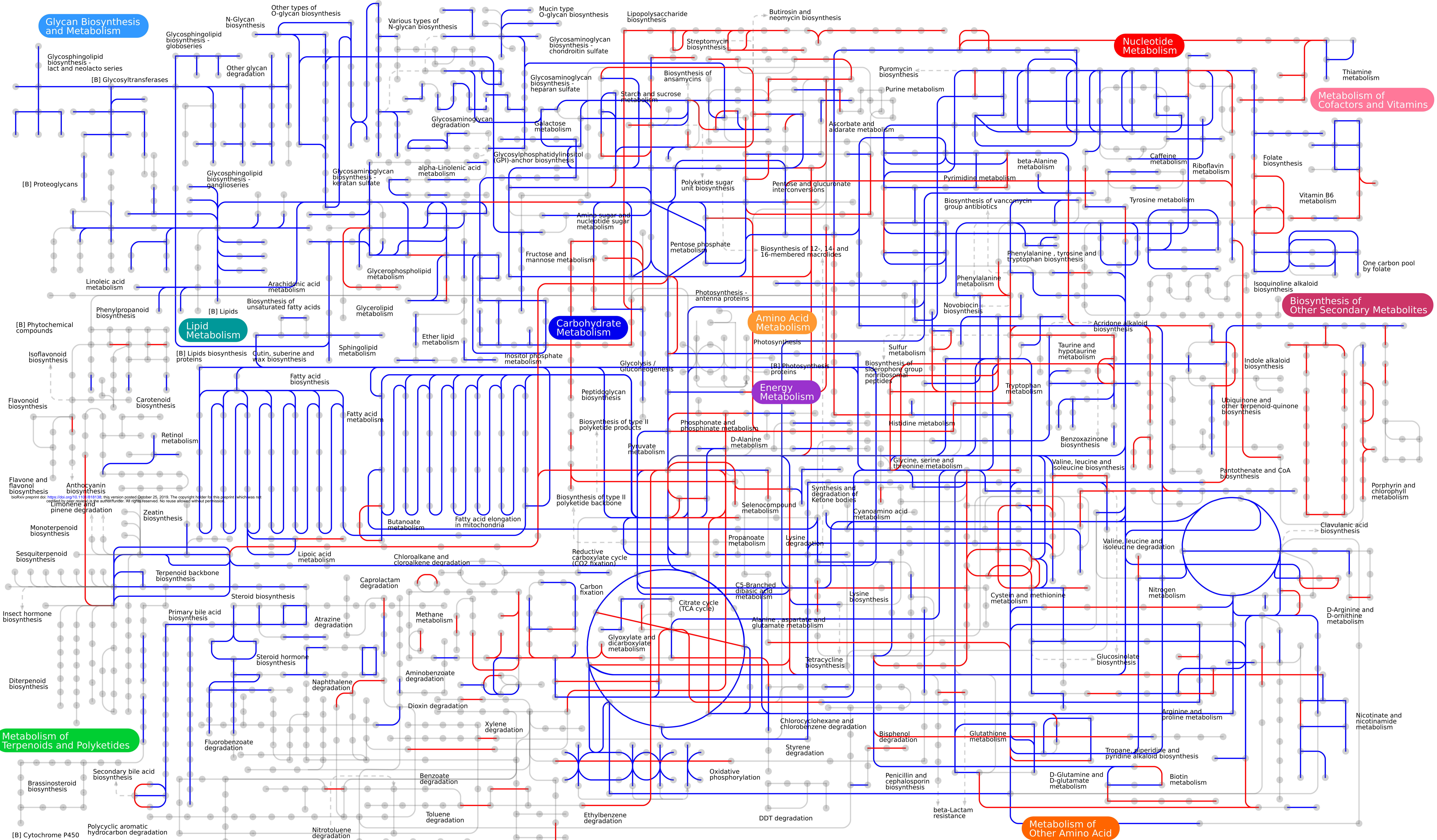
Table 1. Datasets' characteristics

Name	Geographic distribution	Number of individuals	Sequence depth¹	Read length²	Sequencing technology³
<i>Global</i>	Malawi, USA, Venezuela	382	>300K	100	GAllx
<i>Twins</i>	UK	2,727	>15K	2x250	MiSeq
<i>Flemish</i>	Belgium	873	>8K	2x250	MiSeq
<i>Shotgun</i>	USA, Europe, China	123	15M	2x75, 2x100	GAllx, HiSeq

¹ Values represent final sequence depth per sample before analysis (i.e. after chimera removal and subsampling to common depth). ² in bp. ³ Illumina

459





Glycan Biosynthesis and Metabolism

Lipid Metabolism

Carbohydrate Metabolism

Amino Acid Metabolism

Energy Metabolism

Nucleotide Metabolism

Metabolism of Cofactors and Vitamins

Biosynthesis of Other Secondary Metabolites

Metabolism of Other Amino Acid

Metabolism of Terpenoids and Polyketides

Xenobiotics Biodegradation and Metabolism

Metabolism of xenobiotics by cytochrome P450
 Drug metabolism - cytochrome P450
 Drug metabolism - other enzymes

bioRxiv preprint doi: <https://doi.org/10.1101/016136>; this version posted October 25, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.