

1 **A unified sequence catalogue of over 280,000 genomes obtained from the**  
2 **human gut microbiome**

3  
4 Alexandre Almeida<sup>1,2,\*</sup>, Stephen Nayfach<sup>3,4</sup>, Miguel Boland<sup>1</sup>, Francesco Strozzi<sup>5</sup>, Martin  
5 Beracochea<sup>1</sup>, Zhou Jason Shi<sup>6,7</sup>, Katherine S. Pollard<sup>6,7,8,9,10,11</sup>, Donovan H. Parks<sup>12</sup>, Philip  
6 Hugenholtz<sup>12</sup>, Nicola Segata<sup>13</sup>, Nikos C. Kyrpides<sup>3,4</sup> and Robert D. Finn<sup>1,\*</sup>

7  
8 <sup>1</sup>European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, UK;  
9 <sup>2</sup>Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK; <sup>3</sup>U. S. Department of  
10 Energy Joint Genome Institute, Walnut Creek, California, USA; <sup>4</sup>Environmental Genomics and  
11 Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, California,  
12 USA; <sup>5</sup>Enterome Bioscience, Paris, France; <sup>6</sup>Gladstone Institutes, San Francisco, CA, USA;  
13 <sup>7</sup>Chan-Zuckerberg Biohub, San Francisco, CA, USA; <sup>8</sup>Institute for Human Genetics,  
14 University of California San Francisco, San Francisco, CA, USA; <sup>9</sup>Institute for Computational  
15 Health Sciences, University of California San Francisco, San Francisco, CA, USA;  
16 <sup>10</sup>Quantitative Biology Institute, University of California San Francisco, San Francisco, CA,  
17 USA; <sup>11</sup>Department of Epidemiology and Biostatistics, University of California San Francisco,  
18 San Francisco, CA, USA; <sup>12</sup>Australian Centre for Ecogenomics, School of Chemistry and  
19 Molecular Biosciences, The University of Queensland, Queensland, Australia; <sup>13</sup>CIBIO  
20 Department, University of Trento, Trento, Italy.

21  
22 \*Corresponding authors

23 Alexandre Almeida: [aalmeida@ebi.ac.uk](mailto:aalmeida@ebi.ac.uk)

24 Robert D. Finn: [rdf@ebi.ac.uk](mailto:rdf@ebi.ac.uk)

25

## 26 **Abstract**

27 Comprehensive reference data is essential for accurate taxonomic and functional  
28 characterization of the human gut microbiome. Here we present the Unified Human  
29 Gastrointestinal Genome (UHGG) collection, a resource combining 286,997 genomes  
30 representing 4,644 prokaryotic species from the human gut. These genomes contain over 625  
31 million protein sequences used to generate the Unified Human Gastrointestinal Protein  
32 (UHGP) catalogue, a collection that more than doubles the number of gut protein clusters over  
33 the Integrated Gene Catalogue. We find that a large portion of the human gut microbiome  
34 remains to be fully explored, with over 70% of the UHGG species lacking cultured  
35 representatives, and 40% of the UHGP missing meaningful functional annotations. Intra-  
36 species genomic variation analyses revealed a large reservoir of accessory genes and single-  
37 nucleotide variants, many of which were specific to individual human populations. These freely  
38 available genomic resources should greatly facilitate investigations into the human gut  
39 microbiome.

## 40 **Main**

41 The human gut microbiome has been implicated in important phenotypes related to human  
42 health and disease<sup>1,2</sup>. However, incomplete reference data that are missing microbial diversity<sup>3</sup>  
43 hamper our understanding of the roles of individual microbiome species, their interactions and  
44 functions. Hence, establishing a comprehensive collection of microbial reference genomes and  
45 genes is an important step for accurate characterization of the taxonomic and functional  
46 repertoire of the intestinal microbial ecosystem.

47

48 The Human Microbiome Project (HMP)<sup>4</sup> was a pioneering initiative to enrich our knowledge  
49 of human-associated microbiota diversity. Hundreds of genomes from bacterial species with  
50 no sequenced representatives were obtained as part of this project, allowing their use for the  
51 first time in reference-based metagenomic studies. The Integrated Gene Catalogue (IGC)<sup>5</sup> was  
52 subsequently created, combining the sequence data available from the HMP and the  
53 Metagenomics of the Human Intestinal Tract (MetaHIT)<sup>6</sup> consortium. This gene catalogue has  
54 been applied successfully to the study of microbiome associations in different clinical  
55 contexts<sup>7</sup>, revealing microbial composition signatures linked to type 2 diabetes<sup>8</sup>, obesity<sup>9</sup> and  
56 other diseases<sup>10</sup>. But, as the IGC comprises genes with no direct link to their genome of origin,  
57 it lacks contextual data to perform a high-resolution taxonomic classification, establish genetic  
58 linkage and deduce complete functional pathways on a genomic basis.

59

60 Culturing studies have continued to unveil new insights into the biology of our gut  
61 communities<sup>11,12</sup> and are essential for applications in research and biotechnology. However,  
62 the advent of high-throughput sequencing and new metagenomic analysis methods — namely  
63 involving genome assembly and binning — has transformed our understanding of the  
64 microbiome composition both in humans and other environments<sup>13–15</sup>. Metagenomic analyses

65 are able to capture substantial microbial diversity not easily accessible by cultivation by  
66 directly analysing the sample genetic material without the need for culturing, though biases do  
67 exist<sup>16</sup>. Recent studies have massively expanded the known species repertoire of the human  
68 gut, making available unprecedented numbers of new cultured and uncultured genomes<sup>16–20</sup>.  
69 Two culturing efforts isolated and sequenced over 500 human gut-associated bacterial genomes  
70 each<sup>18,20</sup>, while three independent studies<sup>16,17,19</sup> reconstructed 60,000–150,000 microbial  
71 metagenome-assembled genomes (MAGs) from public human microbiome data, most of which  
72 belong to species lacking cultured representatives. Combining these individual efforts and  
73 establishing a unified non-redundant dataset of human gut genomes is essential for driving  
74 future microbiome studies. To accomplish this, we compiled and analysed 286,997 genomes  
75 and 625,251,941 genes from human gut microbiome datasets to generate the Unified Human  
76 Gastrointestinal Genome (UHGG) and Protein (UHGP) catalogues, the most comprehensive  
77 sequence resources of the human gut microbiome established to date.

78

## 79 **Results**

### 80 **The UHGG represents over 280,000 human gut microbial genomes**

81 We first gathered all prokaryotic isolate genomes and MAGs from the human gut microbiome  
82 (publicly available as of March 2019). We compiled the isolate genomes from the Human  
83 Gastrointestinal Bacteria Culture Collection (HBC)<sup>18</sup>, the Culturable Genome Reference  
84 (CGR)<sup>20</sup>, as well as cultured human gut genomes available in the NCBI<sup>21</sup>, PATRIC<sup>22</sup> and  
85 IMG<sup>23</sup> repositories which include genomes from several other large studies<sup>11,12,24</sup>. In addition,  
86 we included all of the gut MAGs generated in Pasolli, et al.<sup>19</sup> (CIBIO), Almeida, et al.<sup>17</sup> (EBI)  
87 and Nayfach, et al.<sup>16</sup> (HGM). To standardize the genome quality across all sets, we used  
88 thresholds of >50% genome completeness and <5% contamination, combined with an  
89 estimated quality score (completeness – 5 × contamination) >50. Final numbers of genomes

90 matching these criteria were: 734 (HBC), 1,519 (CGR), 651 (NCBI), 7,744 (PATRIC/IMG),  
91 137,474 (CIBIO), 87,386 (EBI) and 51,489 (HGM), resulting in a total of 286,997 genome  
92 sequences (Fig. 1a and Supplementary Table 1). Genomes were recovered in samples from a  
93 total of 31 countries across six continents (Africa, Asia, Europe, North America, South  
94 America and Oceania), but the majority originated from samples collected in China, Denmark,  
95 Spain and the United States (Fig. 1b).

96  
97 To determine how many species were included in this gut reference collection, we clustered all  
98 286,997 genomes using a multi-step distance-based approach (see ‘Methods’) with an average  
99 nucleotide identity (ANI) threshold of 95% over at least a 30% alignment fraction<sup>25</sup>. The  
100 clustering procedure resulted in a total of 4,644 inferred prokaryotic species (4,616 bacterial  
101 and 28 archaeal, Supplementary Table 2). We found the species clustering results to be highly  
102 consistent with those previously obtained<sup>16,17,19</sup> (Supplementary Table 3). The best quality  
103 genome from each species cluster was selected as its representative on the basis of genome  
104 completeness, contamination and assembly N50 (with isolate genomes always given  
105 preference over MAGs), and the final set was used to generate the Unified Human  
106 Gastrointestinal Genome (UHGG) catalogue (Fig. 1c). Out of the 4,644 species-level genomes,  
107 3,207 were >90% complete (interquartile range, IQR = 87.2–98.8%) and <5% contaminated  
108 (IQR = 0.0–1.34%), with 573 of those having the 5S, 16S and 23S rRNA genes together with  
109 at least 18 of the standard tRNAs (Supplementary Fig. 1). These 573 genomes satisfy the “high  
110 quality” criteria set for MAGs by the Genomic Standards Consortium<sup>26</sup>. Thereafter, we  
111 classified each species representative using the Genome Taxonomy Database<sup>27</sup> Toolkit  
112 (Supplementary Fig. 2), a standardized taxonomic framework based on a concatenated protein  
113 phylogeny representing >140,000 public prokaryote genomes, fully resolved to the species  
114 level (see ‘Methods’ for details on the taxonomy nomenclature used). However, over 60% of

115 the gut genomes could not be assigned to an existing species, confirming the majority of the  
116 UHGG species lack representation in current reference databases.

117

### 118 **Comparison of species recovered in individual studies**

119 We investigated how many of the 4,644 gut species were found in the different study  
120 collections in order to determine their level of overlap and reproducibility, as well as the ratio  
121 between cultured and uncultured species (Fig. 2a). The largest intersection was between the  
122 collections of MAGs, with the same 1,081 species detected independently in the CIBIO, EBI  
123 and HGM datasets, but not in any of the cultured genome studies. By restricting the analysis to  
124 genomes recovered from 1,554 samples common to all three MAG studies, we found that 93-  
125 97% of species from each set were detected in at least one other MAG collection, and 79-86%  
126 across all three (Supplementary Fig. 3a). Similar level of species overlap was observed when  
127 comparing studies on a per-sample basis (Supplementary Fig. 3b). Further, conspecific  
128 genomes recovered from the same samples across different studies shared a median ANI and  
129 aligned fraction of 99.9% and 92.1%, respectively (Supplementary Fig. 3c). These results  
130 emphasize the reproducibility of the different assembly and binning methods used in the large-  
131 scale studies of human gut MAGs<sup>16,17,19</sup>. Importantly, rarefaction analysis indicates the number  
132 of uncultured species detected has not reached a saturation point, meaning additional species  
133 remain to be discovered (Fig. 2b). However, these most likely represent rarer members of the  
134 human gut microbiome, as the number of species is closer to saturating when only considering  
135 those with at least two conspecific genomes.

136

137 We also investigated the intersection between the three large culture-based datasets: the HBC,  
138 CGR and the NCBI (which contains gut genomes from the Human Microbiome Project,  
139 HMP<sup>4</sup>). Unlike the MAGs, the majority of cultured species were unique within a single

140 collection (486/698; 70%), with only 70 (10%) being common to all three collections  
141 (Supplementary Fig. 3d). This may be due to varied geographical sampling between the  
142 collections (Asia, Europe and North America) or highlight the stochastic nature of culture-  
143 based studies.

144  
145 By calculating the number of genomes contained within each cultured and uncultured species,  
146 we found that species containing isolate genomes represented the largest clusters, while those  
147 exclusively encompassing MAGs tended to be the rarest, as discussed previously<sup>16,17,19</sup>. For  
148 example, only two of the 25 largest bacterial clusters were exclusively represented by MAGs  
149 (Fig. 2c), with 1,212 uncultured species represented by a single genome (80% of which  
150 originated from samples only analysed in one of the MAG studies; Supplementary Fig. 4). The  
151 bacterial species most represented in our collection were *Agathobacter rectalis* (recently  
152 reclassified from *Eubacterium rectale*<sup>28</sup>), *Escherichia coli* D and *Bacteroides uniformis* (Fig.  
153 2c, Supplementary Fig. 5 and Supplementary Table 2), whereas the most frequently recovered  
154 archaeal species was *Methanobrevibacter A smithii*, with 608 genomes found across all six  
155 continents (Supplementary Fig. 6). The largest species clusters displayed similarly high levels  
156 of geographical distribution, indicating the most highly represented species were not restricted  
157 to individual locations (Fig. 2c and Supplementary Fig. 5b).

158

### 159 **Most gut microbial species still lack isolate genomes**

160 We found that 3,750 (81%) of the species in the UHGG did not have a representative in any of  
161 the human gut culture databases. To extend the search to isolate genomes from other  
162 environments or lacking information on their isolation source, we compared the UHGG  
163 catalogue to all NCBI RefSeq isolate genomes. We identified an additional set of 438 species

164 closely matching cultured genomes, leaving 3,312 (71%) of UHGG species as uncultured  
165 (Supplementary Table 2).

166

167 The phylogenetic distribution of the 4,616 bacterial (Fig. 3a) and the 28 archaeal species  
168 (Supplementary Fig. 6) revealed that uncultured species exclusively represented 66% and 31%  
169 of the phylogenetic diversity of Bacteria and Archaea, respectively, with several phyla lacking  
170 cultured representatives (Fig. 3b). The four largest monophyletic groups lacking cultured  
171 genomes were the 4C28d-15 order (167 species, recently proposed as the novel order  
172 Comantemales ord. nov.<sup>29</sup>; Fig. 3c), order RF39 (139 species), family CAG-272 (88 species),  
173 and order Gastranaerophilales (67 species). While none have been successfully cultured,  
174 several have been described in the literature, including RF39<sup>16</sup> and Gastranaerophilales  
175 (previously classified as a lineage in the Melainabacteria<sup>30</sup>) which are characterized by highly  
176 reduced genomes with numerous auxotrophies. This analysis suggests that, despite recent  
177 culture-based studies<sup>11,12,18,20</sup>, much of the diversity in the gut microbiome remains uncultured,  
178 including several large and prevalent clades.

179

### 180 **The UHGP expands the protein universe in the human gut microbiome**

181 Metagenomic approaches have the ability to leverage gene content information not only for  
182 more precise taxonomic analysis, but to also predict the functional capacity of individual  
183 species of interest compared to marker gene-based methods (e.g. relying solely on the 16S  
184 rRNA gene or a limited number of diagnostic genes). We built the Unified Human  
185 Gastrointestinal Protein (UHGP) catalogue with a total of 625,251,941 full-length protein  
186 sequences predicted from the 286,997 genomes here analysed. These were clustered at 50%  
187 (UHGP-50), 90% (UHGP-90), 95% (UHGP-90) and 100% (UHGP-100) amino acid identity,  
188 generating between 5 to 171 million protein clusters (Fig. 1c and Fig. 4a).



189 To determine how comprehensive the UHGP was when compared to existing human gut gene  
190 catalogues, we combined the UHGP-90 ( $n = 13,910,025$  protein clusters) together with the  
191 Integrated Gene Catalogue<sup>5</sup>, a collection of 9.9 million genes from 1,267 gut metagenome  
192 assemblies, which we grouped into 7,063,981 protein clusters at 90% protein identity (referred  
193 to as IGC-90). Nearly all samples used to generate the IGC were also included in the UHGP  
194 catalogue (except for 59 transcriptome datasets), but the latter was generated from a larger and  
195 more geographically diverse metagenomic dataset (including samples from Africa, South  
196 America and Oceania). The UHGP-90 and IGC-90 resulted in a combined set of 15.2 million  
197 protein clusters, with an overlap of 5.8 million sequences (Fig. 4b). This revealed that 81% of  
198 the IGC is represented in the UHGP catalogue, with the missing 19% likely representing  
199 fragments of prokaryotic genomes <50% complete, viral or eukaryotic sequences, plasmids or  
200 other sequences not binned into MAGs. Most notably though, the UHGP provided an increase  
201 of 115% coverage of the gut microbiome protein space over the IGC. As the UHGP was  
202 generated from individual genomes and not from their original unbinned metagenome  
203 assemblies, our catalogue also has the advantage of providing a direct link between each gene  
204 cluster and its genome of origin. This ultimately allows combining individual genes with their  
205 genomic context for an integrated study of the gut microbiome.

206

### 207 **Functional capacity of the human gut microbiota**

208 We used the eggNOG<sup>31</sup>, InterPro<sup>32</sup>, COG<sup>33</sup> and KEGG<sup>34</sup> annotation schemes to capture the full  
209 breadth of functions within the UHGP. However, we found that 42.6% of the UHGP-100 was  
210 poorly characterized, as 28.1% lacked a match to any database and a further 14.5% only had a  
211 match to a COG with no known function (Fig. 4c). Based on the distribution of COG functions,  
212 the most highly represented categories were related to amino acid transport and metabolism,  
213 cell wall/membrane/envelope biogenesis and transcription.

214 We further leveraged the set of 625 million proteins derived from the human gut genomes to  
215 explore the functional diversity within each of the UHGG species. Protein sequences from all  
216 conspecific genomes were clustered at a 90% amino acid identity to generate a pan-genome  
217 for each species. Analysis of the functional capacity of the UHGG species pan-genomes  
218 identified a total of 363 KEGG modules encoded by at least one species (Supplementary Fig.  
219 7a and Supplementary Table 4). Most conserved modules were related to ribosomal structure,  
220 glycolysis, inosine monophosphate biosynthesis, gluconeogenesis, and the shikimate pathway  
221 — all representing essential bacterial functions. However, we found that for certain phyla such  
222 as Myxococcota, Bdellovibrionota, Thermoplasmatota, Patescibacteria and  
223 Verrucomicrobiota, a substantial proportion of the species pan-genomes remained poorly  
224 characterized (Supplementary Fig. 7b). At the same time, species belonging to the clades  
225 Fibrobacterota, Bacteroidota, Firmicutes I, Verrucomicrobiota and Patescibacteria had the  
226 highest proportion of genes encoding carbohydrate-active enzymes (CAZy; Supplementary  
227 Fig. 7b). As most of these lineages are largely represented by uncultured species (Fig. 3b), this  
228 suggests the gut microbiota may harbour many species with important metabolic activities yet  
229 to be cultured and functionally characterized under laboratory conditions.

230

### 231 **Patterns of intra-species genomic diversity**

232 With the protein annotations and pan-genomes inferred for each of the UHGG species, we  
233 explored their intra-species core and accessory gene repertoire. Only near-complete genomes  
234 ( $\geq 90\%$  completeness) and species with at least 10 independent conspecific genomes were  
235 analysed. The overall pattern of gene frequency within each of the 781 species here considered  
236 showed a distinctive bimodal distribution (Supplementary Fig. 8), with most genes classified  
237 as either core or rare (i.e. present in  $\geq 90\%$  or  $< 10\%$  of conspecific genomes). We analysed the  
238 pan-genome size per species in relation to the number of conspecific genomes to look for

239 differences in intra-species gene richness. We observed distinct patterns across different gut  
240 phyla, with species from various Firmicutes clades showing the highest rates of gene gain (Fig.  
241 5a). There was a wide variation in the proportion of core genes between species even among  
242 clades with more than 1,000 genomes (Fig. 5b), with a median core genome proportion  
243 (percentage of core genes out of all genes in the representative genome) estimated at 66% (IQR  
244 = 59.6–73.9%).

245

246 To distinguish the functions encoded in the core and accessory genes, we analysed their  
247 associated annotations. Core genes were well covered, with a median of 96%, 94%, 92% and  
248 69% of the genes assigned with an eggNOG, InterPro, COG and KEGG annotation,  
249 respectively (Fig. 5c). However, the accessory genes had a significantly higher proportion of  
250 unknown functions ( $P < 0.001$ ), with a median of 21% of the genes (IQR = 16.7–27.3%) lacking  
251 a match in any of the databases considered. Thereafter, we investigated the functions encoded  
252 by the core and accessory genes on the basis of the COG functional categories. Genes classified  
253 as core were significantly associated (adjusted  $P < 0.001$ ) with key metabolic functions  
254 involved in nucleotide, amino acid and lipid metabolism, as well as other housekeeping  
255 functions (e.g. related to translation and ribosomal structure, Fig. 5d). In contrast, accessory  
256 genes had a much greater proportion of COGs without a known function, and of genes involved  
257 in replication and recombination which are typically found in mobile genetic elements (MGEs,  
258 Fig. 5d). A significant number of accessory genes were related to defence mechanisms, which  
259 encompass not only general mechanisms of antimicrobial resistance (AMR) such as ABC  
260 transporter efflux pumps, but also targeted systems towards invading MGEs (e.g. CRISPR-Cas  
261 and restriction modification systems against bacteriophages). These results highlight the  
262 potential of this resource to better understand the dynamics of chromosomally encoded AMR

263 within the gut and decipher to what extent the microbiome may be a source of both known and  
264 novel resistance mechanisms.

265

266 We next investigated intra-species single nucleotide variants (SNVs) within the UHGG  
267 species. We generated a catalogue consisting of 249,435,699 SNVs from 2,489 species with  
268 three or more conspecific genomes (Fig. 6a). For context, a previously published catalogue  
269 contained 10.3 million single nucleotide polymorphisms from 101 gut microbiome species<sup>35</sup>.  
270 Of note, more than 85% of these SNVs were exclusively detected in MAGs, whereas only 2.2%  
271 were exclusive to isolate genomes (Fig. 6b). We found the overall pairwise SNV density  
272 between MAGs to be higher than that observed between isolate genomes (Fig. 6c). Next, we  
273 assigned the detected SNVs to the continent of origin of each genome and observed that 36%  
274 of the SNVs were continent exclusive. Notably, genomes with a European origin contributed  
275 to the most exclusive SNVs (Fig. 6d). However, genomes from Africa contributed over three  
276 times more variation on average than European or North American genomes. Pairwise SNV  
277 analysis also supported a higher cross-continent SNV density, especially between genomes  
278 from Africa and Europe (Fig. 6e). Our results suggest there is a high strain variability between  
279 continents and that a considerable level of diversity remains to be discovered, especially from  
280 underrepresented regions such as Africa, South America and Oceania.

281

## 282 **Resource implementation**

283 Both the UHGG and UHGP catalogues are available as part of a new genome layer within the  
284 MGnify<sup>36</sup> website, where summary statistics of each species cluster and their functional  
285 annotations can be interactively explored and downloaded (see ‘Data availability’ section for  
286 more details). We have also generated a BItSliced Genomic Signature Index (BIGSI)<sup>37</sup> of the  
287 UHGG, which will allow users to interactively screen for the presence of small sequence

288 fragments (<5 kb) in this collection. As new genomes from the human gut microbiome are  
289 generated and made publicly available, we plan to periodically update the resource with newly  
290 discovered species or by replacing uncultured reference genomes with better quality versions.

291

## 292 **Discussion**

293 We have generated a unified sequence catalogue representing 286,997 genomes and over 625  
294 million protein sequences of the human gut microbiome. Of the 4,644 species contained in the  
295 UHGG, 71% lack a cultured representative, meaning the majority of microbial diversity in the  
296 catalogue remains to be experimentally characterized. During preparation of our manuscript, a  
297 new collection of almost 4,000 cultured genomes from 106 gut species was released<sup>38</sup>, which  
298 will be incorporated in future versions of the resource. As 96% of these genomes were reported  
299 to have a species representative in the culture collections here included, we do not anticipate  
300 this dataset to provide a substantial increase in the number of species discovered. Nevertheless,  
301 our analyses suggest additional uncultured species from the human gut microbiome are yet to  
302 be discovered, highlighting the importance and need for culture-based studies. Furthermore,  
303 given the sampling bias towards populations from China, Europe and the United States, we  
304 expect that many underrepresented regions still contain substantial uncultured diversity.

305

306 By comparing recently published large datasets of uncultured genomes<sup>16,17,19</sup>, we were able to  
307 assess the reproducibility of the results from each study. We show that despite the different  
308 assembly, binning and refinement procedures employed in the three studies, almost all of the  
309 same species and near-identical strains were recovered independently when using a consistent  
310 sample set. These results further increase confidence in the use of metagenome-assembled  
311 genomes for the characterization of uncultured microbial diversity.

312

313 With the establishment of this massive sequence catalogue, it is evident that a large portion of  
314 the species and functional diversity within the human gut microbiome remains uncharacterized.  
315 Moreover, our knowledge of the intra-species diversity of many species is still limited due to  
316 the presence of a small number of conspecific genomes. Having this combined resource can  
317 help guide future studies and prioritize targets for further experimental validation. Using the  
318 UHGG or UHGP, the community can now screen for the prevalence and abundance of  
319 species/genes in a large panel of intestinal samples and in specific clinical contexts. By  
320 pinpointing particular taxonomic groups with biomedical relevance, more targeted approaches  
321 could be developed to improve our understanding of their role in the human gut. The functional  
322 predictions generated for the species pan-genomes could also be leveraged to develop new  
323 culturing strategies for isolation of candidate species. Target-enrichment methods such as  
324 single-cell<sup>39</sup> and/or bait-capture hybridization<sup>40</sup> approaches could also be applied. Being able  
325 to enrich for specific groups of interest, even without culturing, could allow recovery of better-  
326 quality versions of MAGs and improve the analysis derived from genome sequence data alone.  
327 Given the large uncultured diversity still remaining in the human gut microbiome, having a  
328 high-quality catalogue of all currently known species substantially enhances the resolution and  
329 accuracy of metagenome-based studies. Therefore, the presented genome and protein catalogue  
330 represents a key step towards a hypothesis-driven, mechanistic understanding of the human gut  
331 microbiome.  
332

## 333 **References**

- 334 1. Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes.  
335 *Nature* **490**, 55–60 (2012).
- 336 2. Feng, Q. *et al.* Gut microbiome development along the colorectal adenoma–carcinoma  
337 sequence. *Nat. Commun.* **6**, 6528 (2015).
- 338 3. Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome research.  
339 *BMC Biol.* **17**, 48 (2019).
- 340 4. Human Microbiome Project Consortium, T. Structure, function and diversity of the  
341 healthy human microbiome. *Nature* **486**, 207–214 (2012).
- 342 5. Li, J. *et al.* An integrated catalog of reference genes in the human gut microbiome. *Nat.*  
343 *Biotechnol.* **32**, 834–841 (2014).
- 344 6. Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic  
345 sequencing. *Nature* **464**, 59–65 (2010).
- 346 7. Nayfach, S., Fischbach, M. A. & Pollard, K. S. MetaQuery: a web server for rapid  
347 annotation and quantitative analysis of specific genes in the human gut microbiome.  
348 *Bioinformatics* **31**, 3368–70 (2015).
- 349 8. Wu, H. *et al.* Metformin alters the gut microbiome of individuals with treatment-naive  
350 type 2 diabetes, contributing to the therapeutic effects of the drug. *Nat. Med.* **23**, 850–  
351 858 (2017).
- 352 9. Liu, R. *et al.* Gut microbiome and serum metabolome alterations in obesity and after  
353 weight-loss intervention. *Nat. Med.* **23**, 859–868 (2017).
- 354 10. Armour, C. R., Nayfach, S., Pollard, K. S. & Sharpton, T. J. A metagenomic meta-  
355 analysis reveals functional signatures of health and disease in the human gut  
356 microbiome. *mSystems* **4**, e00332-18 (2019).
- 357 11. Browne, H. P. *et al.* Culturing of ‘unculturable’ human microbiota reveals novel taxa

- 358 and extensive sporulation. *Nature* **533**, 543–546 (2016).
- 359 12. Lagier, J.-C. *et al.* Culture of previously uncultured members of the human gut  
360 microbiota by culturomics. *Nat. Microbiol.* **1**, 16203 (2016).
- 361 13. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes  
362 substantially expands the tree of life. *Nat. Microbiol.* **2**, 1533–1542 (2017).
- 363 14. Stewart, R. D. *et al.* Compendium of 4,941 rumen metagenome-assembled genomes for  
364 rumen microbiome biology and enzyme discovery. *Nat. Biotechnol.* **37**, 953–961  
365 (2019).
- 366 15. Anantharaman, K. *et al.* Thousands of microbial genomes shed light on interconnected  
367 biogeochemical processes in an aquifer system. *Nat. Commun.* **7**, 13219 (2016).
- 368 16. Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. New insights from  
369 uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510  
370 (2019).
- 371 17. Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**,  
372 499–504 (2019).
- 373 18. Forster, S. C. *et al.* A human gut bacterial genome and culture collection for improved  
374 metagenomic analyses. *Nat. Biotechnol.* **37**, 186–192 (2019).
- 375 19. Pasolli, E. *et al.* Extensive unexplored human microbiome diversity revealed by over  
376 150,000 genomes from metagenomes spanning age, geography, and lifestyle. *Cell* **176**,  
377 649–662.e20 (2019).
- 378 20. Zou, Y. *et al.* 1,520 reference genomes from cultivated human gut bacteria enable  
379 functional microbiome analyses. *Nat. Biotechnol.* **37**, 179–185 (2019).
- 380 21. Kitts, P. A. *et al.* Assembly: a resource for assembled genomes at NCBI. *Nucleic Acids*  
381 *Res.* **44**, D73–D80 (2016).
- 382 22. Wattam, A. R. *et al.* Improvements to PATRIC, the all-bacterial bioinformatics database



- 383 and analysis resource center. *Nucleic Acids Res.* **45**, D535–D542 (2017).
- 384 23. Chen, I.-M. A. *et al.* IMG/M v.5.0: an integrated data management and comparative  
385 analysis system for microbial genomes and microbiomes. *Nucleic Acids Res.* **47**, D666–  
386 D677 (2019).
- 387 24. Human Microbiome Jumpstart Reference Strains Consortium, T. H. M. J. R. S. *et al.* A  
388 catalog of reference genomes from the human microbiome. *Science* **328**, 994–9 (2010).
- 389 25. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High  
390 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.  
391 *Nat. Commun.* **9**, 5114 (2018).
- 392 26. Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG)  
393 and a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat.*  
394 *Biotechnol.* **35**, 725–731 (2017).
- 395 27. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny  
396 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996 (2018).
- 397 28. Rosero, J. A. *et al.* Reclassification of *Eubacterium rectale* (Hauduroy *et al.* 1937)  
398 Prévot 1938 in a new genus *Agathobacter* gen. nov. as *Agathobacter rectalis* comb.  
399 nov., and description of *Agathobacter ruminis* sp. nov., isolated from the rumen con.  
400 *Int. J. Syst. Evol. Microbiol.* **66**, 768–773 (2016).
- 401 29. Hildebrand, F. *et al.* Antibiotics-induced monodominance of a novel gut bacterial order.  
402 *Gut* (2019). doi:10.1136/gutjnl-2018-317715
- 403 30. Di Rienzi, S. C. *et al.* The human gut and groundwater harbor non-photosynthetic  
404 bacteria belonging to a new candidate phylum sibling to Cyanobacteria. *Elife* **2**, (2013).
- 405 31. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically  
406 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*  
407 *Res.* **47**, D309–D314 (2019).

- 408 32. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.  
409 *Bioinformatics* **30**, 1236–1240 (2014).
- 410 33. Galperin, M. Y., Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Expanded microbial  
411 genome coverage and improved protein family annotation in the COG database. *Nucleic*  
412 *Acids Res.* **43**, D261–D269 (2015).
- 413 34. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new  
414 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–  
415 D361 (2017).
- 416 35. Schloissnig, S. *et al.* Genomic variation landscape of the human gut microbiome. *Nature*  
417 **493**, 45–50 (2013).
- 418 36. Mitchell, A. L. *et al.* EBI Metagenomics in 2017: enriching the analysis of microbial  
419 communities, from sequence reads to assemblies. *Nucleic Acids Res.* **46**, D726–D735  
420 (2017).
- 421 37. Bradley, P., den Bakker, H. C., Rocha, E. P. C., McVean, G. & Iqbal, Z. Ultrafast search  
422 of all deposited bacterial and viral genomic data. *Nat. Biotechnol.* **37**, 152–159 (2019).
- 423 38. Poyet, M. *et al.* A library of human gut bacterial isolates paired with longitudinal  
424 multiomics data enables mechanistic microbiome research. *Nat. Med.* 1–11 (2019).  
425 doi:10.1038/s41591-019-0559-3
- 426 39. Xu, Y. & Zhao, F. Single-cell metagenomics: challenges and applications. *Protein Cell*  
427 **9**, 501–510 (2018).
- 428 40. Noyes, N. R. *et al.* Enrichment allows identification of diverse, rare elements in  
429 metagenomic resistome-virulome sequencing. *Microbiome* **5**, 142 (2017).
- 430 41. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
431 assessing the quality of microbial genomes recovered from isolates, single cells, and  
432 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).

- 433 42. Nawrocki, E. P., Kolbe, D. L. & Eddy, S. R. Infernal 1.0: inference of RNA alignments.  
434 *Bioinformatics* **25**, 1335–1337 (2009).
- 435 43. Kalvari, I. *et al.* Rfam 13.0: shifting to a genome-centric resource for non-coding RNA  
436 families. *Nucleic Acids Res.* **46**, D335–D342 (2018).
- 437 44. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer  
438 RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–64 (1997).
- 439 45. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate  
440 genomic comparisons that enables improved genome recovery from metagenomes  
441 through de-replication. *ISME J.* **11**, 2864–2868 (2017).
- 442 46. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,  
443 2825–2830 (2011).
- 444 47. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.*  
445 **5**, R12 (2004).
- 446 48. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
447 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 448 49. Turner, I., Garimella, K. V, Iqbal, Z. & McVean, G. Integrating long-range connectivity  
449 information into de Bruijn graphs. *Bioinformatics* **34**, 2556–2565 (2018).
- 450 50. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–  
451 2069 (2014).
- 452 51. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site  
453 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 454 52. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.  
455 *Bioinformatics* **31**, 3691–3693 (2015).
- 456 53. Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time. *Nat.*  
457 *Commun.* **9**, 2542 (2018).

- 458 54. Huerta-Cepas, J. *et al.* Fast genome-wide functional annotation through orthology  
459 assignment by eggNOG-mapper. *Mol. Biol. Evol.* **34**, 2115–2122 (2017).
- 460 55. Lombard, V., Golaconda Ramulu, H., Drula, E., Coutinho, P. M. & Henrissat, B. The  
461 carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res.* **42**, D490–  
462 D495 (2014).
- 463 56. Torchiano, M. Effsize - a package for efficient effect size computation. (2016).  
464 doi:10.5281/ZENODO.1480624
- 465 57. Gloor, G. B., Wu, J. R., Pawlowsky-Glahn, V. & Egozcue, J. J. It's all relative: analyzing  
466 microbiome data as compositions. *Ann. Epidemiol.* **26**, 322–329 (2016).
- 467

## 468 **Methods**

### 469 **Genome collection**

470 We compiled all the prokaryotic genomes publicly available as of March 2019 that have been  
471 sampled from the human gut. To retrieve isolate genomes, we surveyed the IMG<sup>23</sup>, NCBI<sup>21</sup>  
472 and PATRIC<sup>22</sup> databases for genome sequences annotated as having been isolated from the  
473 human gastrointestinal tract. We complemented this set with bacterial genomes belonging to  
474 two recent culturomics collections: the Human Gastrointestinal Bacteria Culture Collection  
475 (HBC)<sup>18</sup> and the Culturable Genome Reference (CGR)<sup>20</sup>. To avoid including duplicated entries  
476 due to redundancy between reference databases, we combined genomes obtained from the  
477 PATRIC and IMG repositories, and added only those without an identical genome in the sets  
478 extracted from NCBI, HBC and CGR. Metagenome-assembled genomes (MAGs, i.e.  
479 uncultured genomes) were obtained from Pasolli, et al.<sup>19</sup> (CIBIO), Almeida, et al.<sup>17</sup> (EBI) and  
480 Nayfach, et al.<sup>16</sup> (HGM). For the CIBIO set, only those genomes retrieved from samples  
481 collected from the intestinal tract were used. Metadata for each genome was retrieved using  
482 the API of the various public repositories and combined with that available in each of the  
483 original studies.

484

### 485 **Assessing genome quality**

486 Genome quality (completeness and contamination) was estimated with CheckM v1.0.11<sup>41</sup>  
487 using the 'lineage\_wf' workflow to select only those that passed the following criteria: >50%  
488 genome completeness, <5% contamination and an estimated quality score (completeness – 5 ×  
489 contamination) >50. We also searched for the presence of ribosomal RNAs in each genome  
490 with the 'cmsearch' function of INFERNAL<sup>42</sup> (options '-Z 1000 --hmmonly --cut\_ga --noali –  
491 tblout') against the Rfam<sup>43</sup> covariance models for the 5S, 16S and 23S rRNAs. tRNAs of the

492 standard 20 amino acids were identified with tRNAScan-SE<sup>44</sup> with options ‘-A -Q’ for archaeal  
493 species and ‘-B -Q’ for those belonging to bacterial lineages.

494

### 495 **Species clustering**

496 We clustered the total set of 286,997 genomes at an estimated species level (average nucleotide  
497 identity, ANI  $\geq 95\%$ <sup>25</sup>) using dRep v2.2.4<sup>45</sup> with the following options: ‘-pa 0.9 -sa 0.95 -nc  
498 0.30 -cm larger’. Because of the computational burden of clustering together the entire genome  
499 set, we employed an iterative approach where random chunks of 50,000 genomes were  
500 clustered independently. The selected representatives from each chunk were combined and  
501 subsequently clustered, reducing the final computational load. To ensure the best quality  
502 genome was selected as the species representative in each iteration, a score was calculated for  
503 each genome based on the following formula:

$$504 \text{ Score} = \text{CMP} - 5 \times \text{CNT} + 0.5 \times \log(\text{N50})$$

505 where CMP represents the completeness level, CNT the estimated contamination and N50 the  
506 assembly contiguity characterized by the minimum contig size in which half of the total  
507 genome sequence is contained. The genome with the highest score was chosen as the species  
508 representative, with cultured genomes prioritized over uncultured genomes (i.e. if a MAG had  
509 a higher score than an isolate genome, the latter would still be chosen as the representative).

510

### 511 **Evaluating methods reproducibility**

512 The species clusters inferred here were compared with those previously generated in the human  
513 gut MAG studies<sup>16,17,19</sup> from a common set of genomes. Similarity between species clusterings  
514 was estimated using the Adjusted Rand Index (ARI) computed in the Scikit-learn python  
515 package<sup>46</sup>. This metric considers both the number of clusters and cluster membership to  
516 compute a similarity score ranging from 0 to 1.

517 Conspecific genomes recovered in the same metagenomic samples but in different studies were  
518 compared with FastANI v1.1<sup>25</sup> with default parameters to obtain both the maximum aligned  
519 fraction and ANI for each pairwise comparison.

520

### 521 **Inferring cultured status**

522 To determine their cultured status, the UHGG species representatives were searched against  
523 NCBI RefSeq release 93 after excluding uncultured genomes (i.e. metagenome-assembled or  
524 single-cell amplified genomes). Genome alignments were performed in two stages: (1) Mash  
525 v2.1 was used as an initial screen (using the function ‘mash dist’) to identify the most similar  
526 RefSeq genome to each of the UHGG species, and (2) ‘dnadiff’ from MUMmer v4.0.0beta2<sup>47</sup>  
527 was subsequently used to compute whole genome ANI between the genome pairs. A species  
528 was considered to have been cultured if (1) it contained a cultured gut genome from the UHGG  
529 catalogue, or (2) if it matched an isolate RefSeq genome with at least 95% ANI over at least  
530 30% of the genome length.

531

### 532 **Calculating number of conspecific genomes**

533 For an accurate assessment of the number of non-redundant genomes belonging to each  
534 species, we de-replicated all conspecific genomes at a 99.9% ANI threshold using dRep with  
535 options ‘-pa 0.999 –SkipSecondary’. Furthermore, the frequency of each species was only  
536 counted once per sample to avoid cases where the same genome was recovered multiple times  
537 because of overlapping samples between the three MAG studies.

538

### 539 **Estimating geographical diversity**

540 A geographical diversity index was estimated to assess how widely distributed each species  
541 was. We calculated the Shannon diversity index on the proportion of samples each species was

542 found per continent. This metric combines both richness and evenness, so the level of estimated  
543 diversity is highest in species found across all continents at a similar proportion.

544

### 545 **Phylogenetic analyses**

546 Taxonomic annotation of each species representative was performed with the Genome

547 Taxonomy Database<sup>27</sup> Toolkit (GTDB-Tk) v0.3.1

548 (<https://github.com/Ecogenomics/GTDBTk>) (database release 04-RS89) using the

549 ‘classify\_wf’ function and default parameters. To use consistent species boundaries between

550 the genome clustering and taxonomic classification procedures, genomes were assigned at the

551 species level if the ANI to the closest GTDB-Tk species representative genome was  $\geq 95\%$  and

552 the alignment fraction  $\geq 30\%$ . In this taxonomy scheme, genera and species names with an

553 alphabetic suffix indicate taxa that are polyphyletic or needed to be subdivided based on

554 taxonomic rank normalization according to the current GTDB reference tree. The lineage

555 containing the type strain retains the unaffixed (valid) name and all other lineages are given

556 alphabetic suffixes, indicating they are placeholder names that need to be replaced in due

557 course. Taxon names above the rank of genus appended with an alphabetic suffix indicate

558 groups that are not monophyletic in the GTDB reference tree, but for which there exists

559 alternative evidence that they are monophyletic groups. We also generated NCBI taxonomy

560 annotations for each species-level genome based on its placement in the GTDB tree, using the

561 ‘gtdb\_to\_ncbi\_majority\_vote.py’ script available in the GTDB-Tk repository

562 (<https://github.com/Ecogenomics/GTDBTk/tree/stable/scripts>).

563

564 Maximum-likelihood trees were generated *de novo* using the protein sequence alignments

565 produced by the GTDB-Tk: we used IQ-TREE v1.6.11 to build a phylogenetic tree of the 4,616

566 bacterial and 28 archaeal species. The best-fit model was automatically selected by



567 ‘ModelFinder’ on the basis of the Bayesian Information Criterion (BIC) score. The LG+F+R10  
568 model was chosen for building the bacterial tree, while the LG+F+R4 model was used for the  
569 archaeal phylogeny. Trees were visualized and annotated with the Interactive Tree Of Life  
570 (iTOL) v4.4.2<sup>48</sup>. Phylogenetic diversity (PD) was estimated by the sum of branch lengths, with  
571 the amount that was exclusive to uncultured species calculated as  $PD_{total} - PD_{cultured}$ . Uncultured  
572 monophyletic groups were defined as nodes in the tree containing child leaves exclusively  
573 comprised of uncultured genomes.

574

### 575 **BIGSI construction**

576 A Bitsliced Genomic Signature Index (BIGSI)<sup>37</sup> was generated for all species-level genomes  
577 with BIGSI v0.3.8. First, *k*-mers of size 31 were extracted from each genome with McCortex  
578 v1.0.1<sup>49</sup> (‘mccortex31 build -k 31’). Thereafter, Bloom filters were built for each *k*-mer set  
579 using ‘bigsi bloom’ and inserted into the BIGSI index with ‘bigsi build’. BIGSI config  
580 parameters *h* (number of hash functions applied to each *k*-mer) and *m* (Bloom filter’s length in  
581 bits) were set at 1 and 28,000,000, respectively. A final API layer for querying the index was  
582 built using hug (<http://www.hug.rest/>) and hosted on the MGnify<sup>36</sup> website:  
583 <https://www.ebi.ac.uk/metagenomics/genomes>.

584

### 585 **Pan-genome analysis and functional annotation**

586 Protein coding sequences (CDS) for each of the 286,997 genomes were predicted and annotated  
587 with Prokka v1.13.3<sup>50</sup>, using Prodigal v2.6.3<sup>51</sup> with options ‘-c’ (predict proteins with closed  
588 ends only), ‘-m’ (prevent genes from being built across stretches of sequences marked as Ns)  
589 and ‘-p single’ (single mode for genome assemblies containing one single species). Pan-  
590 genome analyses were carried out using Roary v3.12.0<sup>52</sup>. We set a minimum amino acid  
591 identity for a positive match at 90% (‘-i 90’), a core gene defined at 90% presence (‘-cd 90’)

592 and no paralog splitting ('-s'). A normalized pan-genome size was estimated by dividing the  
593 total number of core and accessory genes by the number of genes contained in the species  
594 representative genome.

595

596 The Unified Human Gastrointestinal Protein (UHGP) catalogue was generated from the  
597 combined set of 625,251,941 CDS predicted. Protein clustering of the UHGP and the Integrated  
598 Gene Catalogue (IGC)<sup>5</sup> was performed with the 'linclust' function of MMseqs2 v6-f5a1c<sup>53</sup>  
599 with options: '--cov-mode 1 -c 0.8' (minimum coverage threshold of 80% length of the shortest  
600 sequence) and '--kmer-per-seq 80' (number of *k*-mers selected per sequence, increased from  
601 the default of 21 to improve clustering sensitivity). The '--min-seq-id' option was set at 1, 0.95,  
602 0.9 and 0.5 to generate the catalogues at 100%, 95%, 90% and 50% protein identity,  
603 respectively. We clustered the IGC solely at a 90% and 50% protein identity as it was originally  
604 de-replicated at a 95% nucleotide identity<sup>5</sup>. Functional characterization of all protein sequences  
605 was performed with eggNOG-mapper v2<sup>54</sup> (database v5.0<sup>31</sup>) and InterProScan v5.35-74.0<sup>32</sup>.  
606 COG<sup>33</sup>, KEGG<sup>34</sup> and CAZy<sup>55</sup> annotations were derived from the eggNOG-mapper results.  
607 Differences in annotation coverage and COG functional categories between the core and  
608 accessory genes were evaluated with a two-tailed Wilcoxon rank-sum test in R v3.6.0 (function  
609 'wilcox.test'). Expected *P* values were corrected for multiple testing with the Benjamini-  
610 Hochberg method. Cohen's *d* effect sizes were estimated with the function 'cohen.d' from the  
611 Effsize<sup>56</sup> R package. To accurately estimate the proportion of each KEGG module in the  
612 species pan-genome, we used the compositional data analysis R package CoDaSeq<sup>57</sup>. Pseudo  
613 counts for zero-count data were first imputed using a Bayesian-Multiplicative simple  
614 replacement procedure implemented in the 'cmultRepl' function (method 'CZM'). Final counts  
615 were thereby converted to centred log-ratios using the 'codaSeq.clr' function to account for the  
616 compositional nature of the data and for differences in pan-genome size.

## 617 **SNV analyses**

618 A total of 2,489 species with at least three conspecific genomes were used for generating a  
619 catalogue of single nucleotide variants (SNVs). For each species, we mapped all conspecific  
620 genomes to the representative genome using the ‘nucmer’ program from MUMmer  
621 v4.0.0.beta2<sup>47</sup> and filtered alignments using the ‘delta-filter’ program with options ‘-q -r’ to  
622 exclude chance- and repeat-induced alignments. Thereafter, we identified SNVs using the  
623 ‘show-snps’ program. Single base insertions and deletions were not counted as SNVs. Each  
624 SNV locus was included in the catalogue only when the alternate allele was detected in at least  
625 two conspecific genomes. The final SNV catalogue was generated by unifying the SNV  
626 coordinates on the basis of their position in the species representative genome. The SNV entries  
627 in the catalogue were characterized as genome type-specific or continent-specific based on  
628 whether the alternate allele could be found solely in genomes from a specific genome type or  
629 continent. The number of continent-specific SNVs was normalized by the number of genomes  
630 from the corresponding continent to estimate the contribution per genome to the continent-  
631 specific SNV discoveries.

632  
633 Similar programs and parameters were used for the pairwise genome alignment, but in this case  
634 only near-complete genomes ( $\geq 90\%$  completeness) and species with at least 10 independent  
635 conspecific genomes were considered. Due to the high computational demand, pairwise  
636 alignments of species encompassing more than 1,000 genomes were limited to the best-quality  
637 1,000 genomes. A total of 29,283,684 pairwise genome alignments were performed between  
638 almost 113,000 genomes from 909 species. For each pairwise comparison, we estimated the  
639 total number of SNVs and the overall density as the number of SNVs per kb. In addition, the  
640 pairwise comparisons were organized based on the type and the continent origin of the genomes  
641 in the pair for further downstream analyses. A two-tailed Wilcoxon rank-sum test was used to

642 evaluate differences in SNV distributions. Resulting  $P$  values were corrected for multiple  
643 testing with the Benjamini-Hochberg method.

644

#### 645 **Data availability**

646 Genome assemblies of the UHGG have been deposited in the European Nucleotide Archive  
647 under study accession ERP116715. The UHGG, UHGP and SNV catalogues are available in a  
648 public FTP server ([http://ftp.ebi.ac.uk/pub/databases/metagenomics/mgnify\\_genomes/](http://ftp.ebi.ac.uk/pub/databases/metagenomics/mgnify_genomes/))  
649 alongside functional annotations and the pan-genome results. These data together with the  
650 BIGSI search index of the UHGG can also be accessed interactively on the MGnify website:  
651 <https://www.ebi.ac.uk/metagenomics/genomes>.

652

#### 653 **Acknowledgements**

654 We would like to thank Phelim Bradley and Zamin Iqbal for their help in the BIGSI  
655 implementation, and Dongying Wu for assistance in the identification of uncultured  
656 monophyletic groups. Funding: European Molecular Biology Laboratory (EMBL);  
657 Biotechnology and Biological Sciences Research Council [BB/N018354/1 and  
658 BB/R015228/1]; European Research Council (project ERC- STG MetaPG-716575) to N.S.

659

#### 660 **Author contributions**

661 A.A., S.N., N.K.C. and R.D.F. conceived the study. A.A. performed the genome clustering and  
662 annotations, compared study sets, carried out the pan-genome analyses, built the BIGSI index  
663 and drafted the manuscript. S.N. provided feedback, performed phylogenetic, rarefaction and  
664 clustering analyses, as well as the comparison with RefSeq. M.Boland and M.Beracochea built  
665 the resource implementation within the MGnify website. F.S. built the protein catalogue and  
666 performed the comparison with the IGC. Z.J.S. generated the SNV catalogue and performed

667 related analyses according to genome type and geographic origin. K.S.P. provided feedback,  
668 funding and contributed to the SNV analyses. D.H.P. and P.H. provided feedback and assisted  
669 in the species taxonomic classification. N.S. provided feedback, funding and contributed to the  
670 generation of the protein catalogue. N.K.C. and R.D.F. supervised the work, provided feedback  
671 and funding. All authors read, edited and approved the final manuscript.

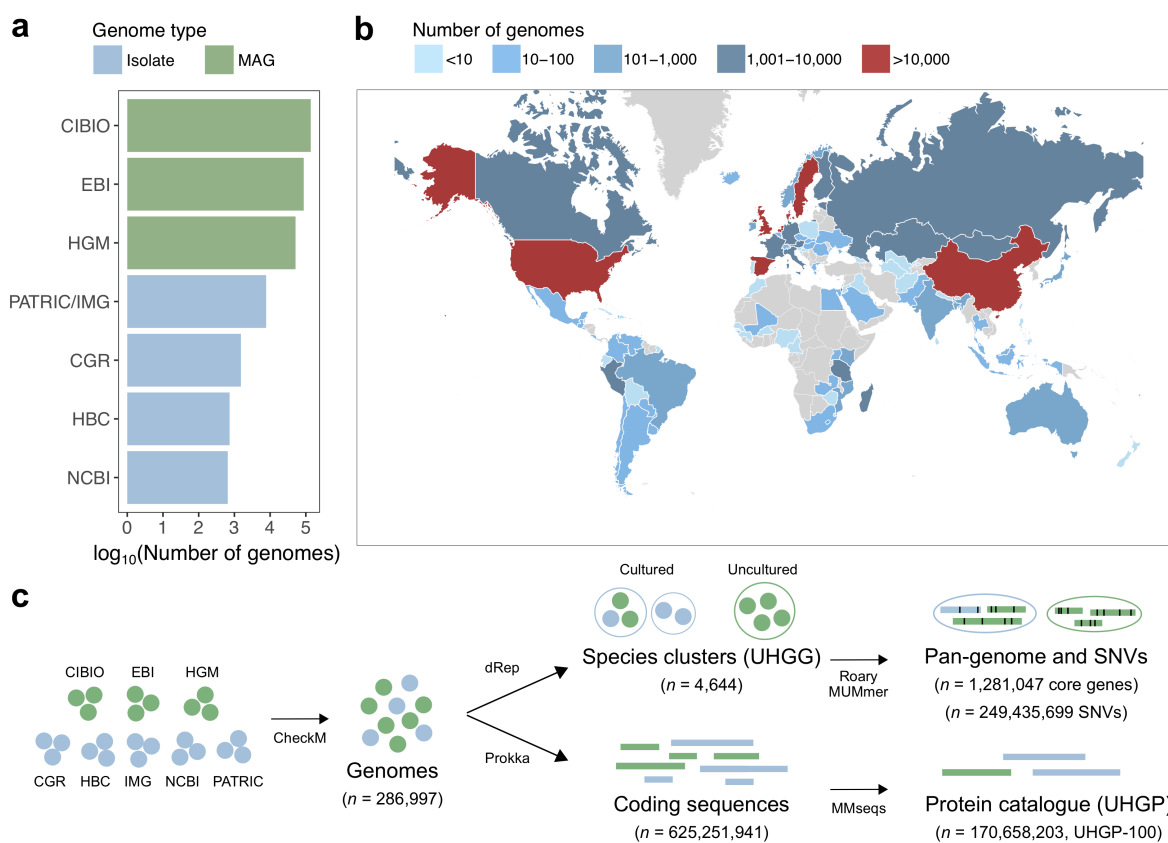
672

### 673 **Competing interests**

674 F.S. is an employee of Enterome SA. P.H. is a co-founder and Director of Microba Life  
675 Sciences Limited. D.H.P. is a consultant to Microba Life Sciences Limited. R.D.F. is a  
676 consultant to Microbiotica Pty Ltd.

677

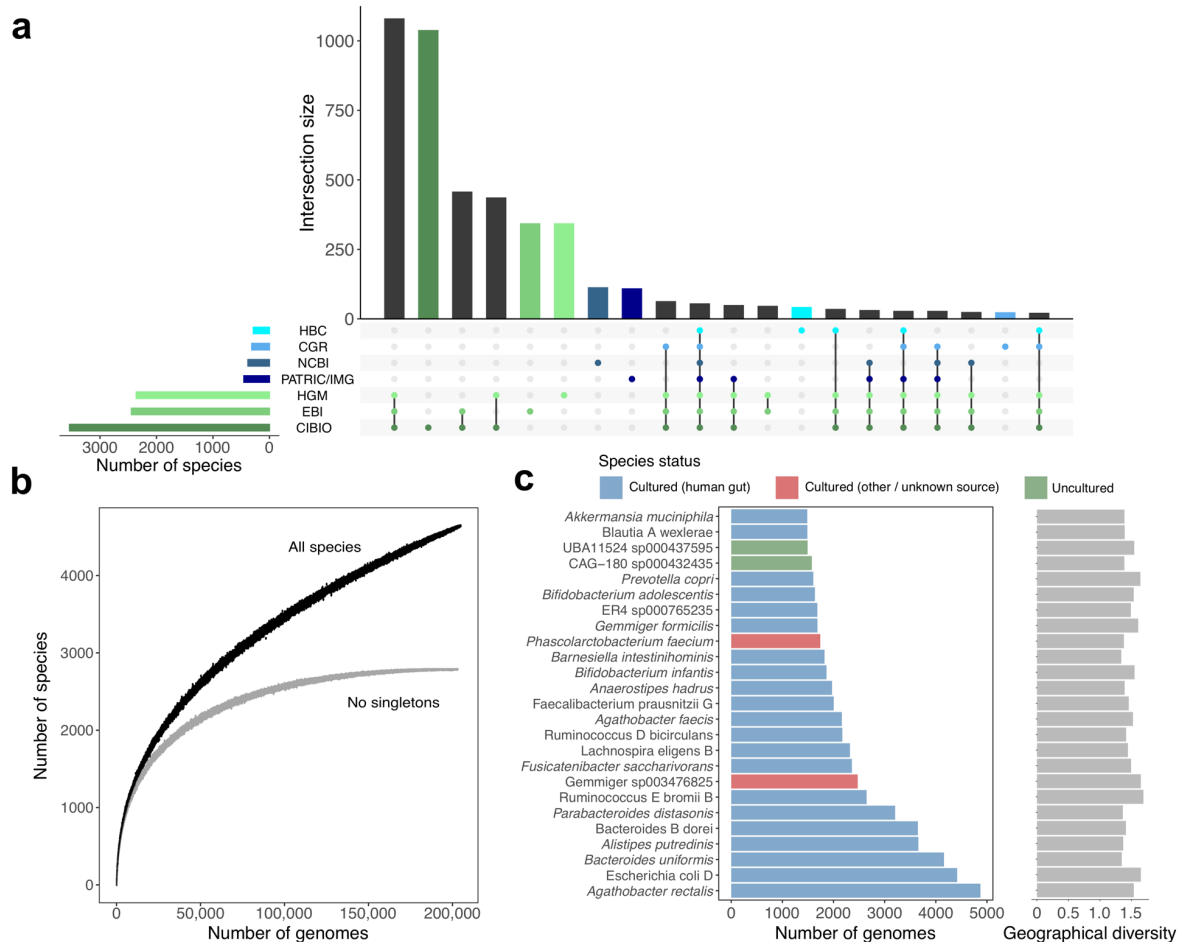
## 678 Figures



679

### 680 **Figure 1. The unified sequence catalogue of the human gut microbiome.**

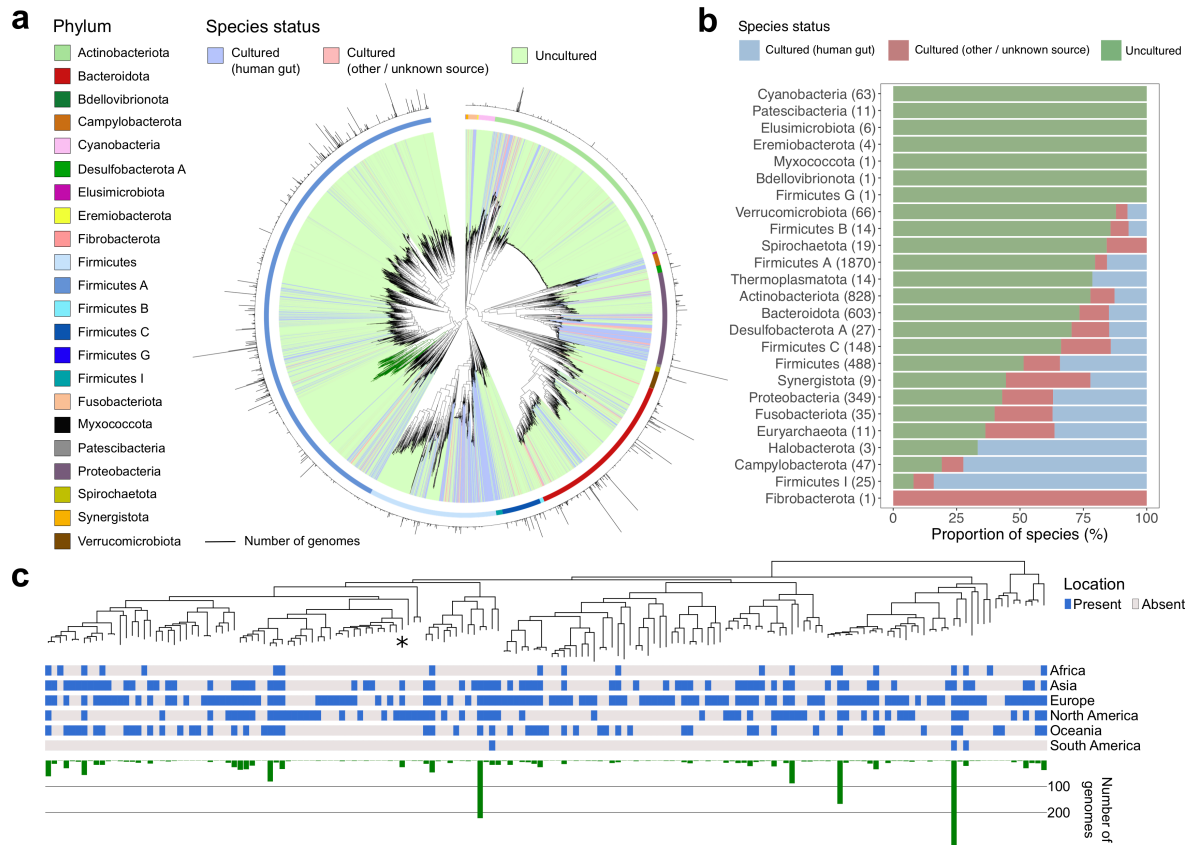
681 **a**, Number of gut genomes per each study set used to generate the sequence catalogues,  
 682 coloured according to whether they represent isolate genomes or metagenome-assembled  
 683 genomes (MAGs). **b**, Geographic distribution of the number of genomes retrieved per country.  
 684 **c**, Overview of the methods used to generate the genome (UHGG) and protein sequence  
 685 (UHGP) catalogue. Genomes retrieved from public datasets were first quality-controlled by  
 686 CheckM. Filtered genomes were clustered at an estimated species-level (95% average  
 687 nucleotide identity) and their intra-species diversity was assessed (genes from conspecific  
 688 genomes were clustered at a 90% protein identity). In parallel, a non-redundant protein  
 689 catalogue was generated from all the coding sequences of the 286,997 genomes at 100%  
 690 (UHGP-100,  $n = 170,658,203$ ), 95% (UHGP-95,  $n = 20,240,320$ ), 90% (UHGP-90,  $n =$   
 691  $13,910,025$ ) and 50% (UHGP-50,  $n = 4,736,012$ ) protein identity.



692

693 **Figure 2. Intersection and frequency of species across studies.**

694 **a**, Number of species found across the genome study sets here used, ordered by their level of  
 695 overlap. Vertical bars represent the number of species shared between the study sets  
 696 highlighted in the lower panel. **b**, Rarefaction curve of the number of species detected as a  
 697 function of the number of non-redundant genomes analysed. Curves are depicted both for all  
 698 the UHGG species, and after excluding singleton species (represented by only one genome). **c**,  
 699 Number of non-redundant genomes detected per species (left) alongside the degree of  
 700 geographical diversity (calculated with the Shannon diversity index, right).



701

702 **Figure 3. Uncultured species are predominant among human gut phyla.**

703 **a**, Maximum-likelihood phylogenetic tree of the 4,616 bacterial species detected in the human

704 gut. Clades are coloured by species cultured status with outer circles depicting the GTDB

705 phylum annotation. Bar graphs in the outermost layer indicate the number of genomes from

706 each species. The order Comantemales ord. nov. is highlighted with dark green branches. **b**,

707 Proportion of species within the 25 prokaryotic phyla detected according to their cultured

708 status. Numbers in brackets represent the total number of species in the corresponding phylum.

709 **c**, Phylogenetic tree of species belonging to the order Comantemales ord. nov. (phylum

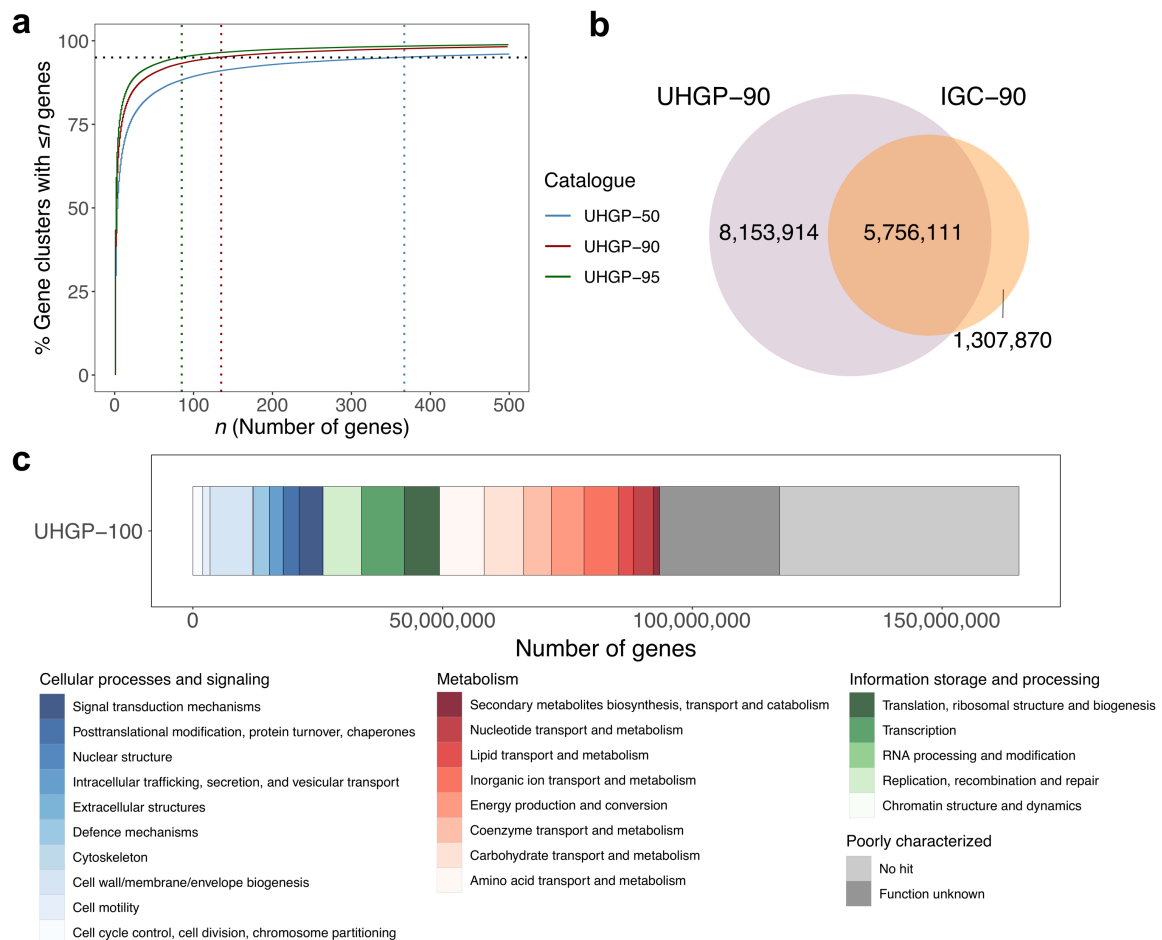
710 Firmicutes A), the largest phylogenetic group exclusively represented by uncultured species.

711 The geographic distribution of each species and the number of genomes recovered is

712 represented below the tree. The species previously classified as *Candidatus* Borkfalki

713 *ceftriaxensis* is indicated with an asterisk.

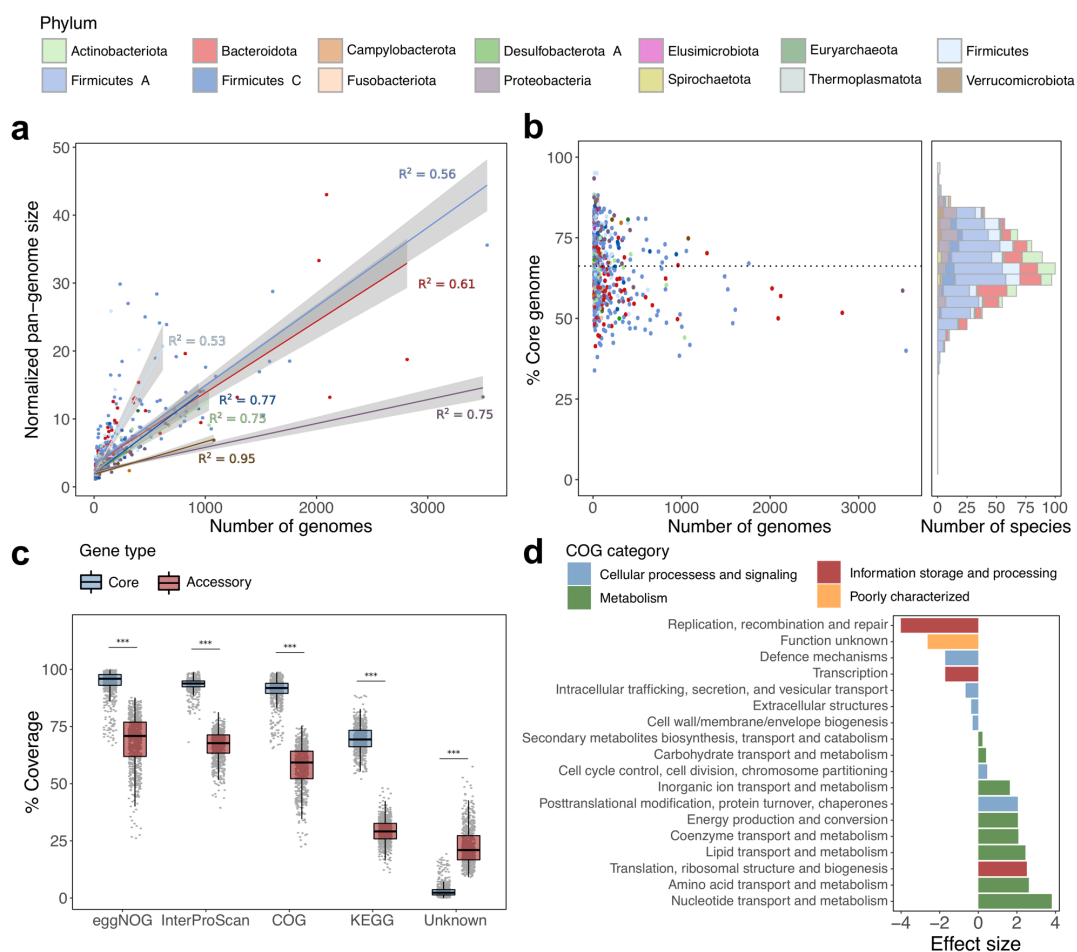




714

715 **Figure 4. The UHGP improves coverage of the human gut protein landscape.**

716 **a**, Cumulative distribution curve of the number and size of the gene clusters of the UHGP-95  
 717 ( $n = 20,240,320$ ), UHGP-90 ( $n = 13,910,025$ ) and UHGP-50 ( $n = 4,736,012$ ). Dashed vertical  
 718 lines indicate the cluster size below which 90% of the gene clusters can be found. **b**, Overlap  
 719 between the UHGP (purple) and IGC (orange), both clustered at 90% amino acid identity. **c**,  
 720 COG functional annotation results of the unified gastrointestinal protein catalogue clustered at  
 721 100% amino acid identity (UHGP-100).



722

723 **Figure 5. Pan-genome diversity patterns within the gut microbiome.**

724 **a**, Normalized pan-genome size as a function of the number of conspecific genomes.

725 Regression curves were generated per phylum, with the corresponding coefficients of

726 determination indicated next to each curve. **b**, Fraction of each species core genome

727 (proportion of core genes out of all genes in the representative genome) according to the

728 number of conspecific genomes (left) and as a histogram (right), coloured by phylum.

729 Horizontal dashed line represents the median value across all species. **c**, Proportion of core and

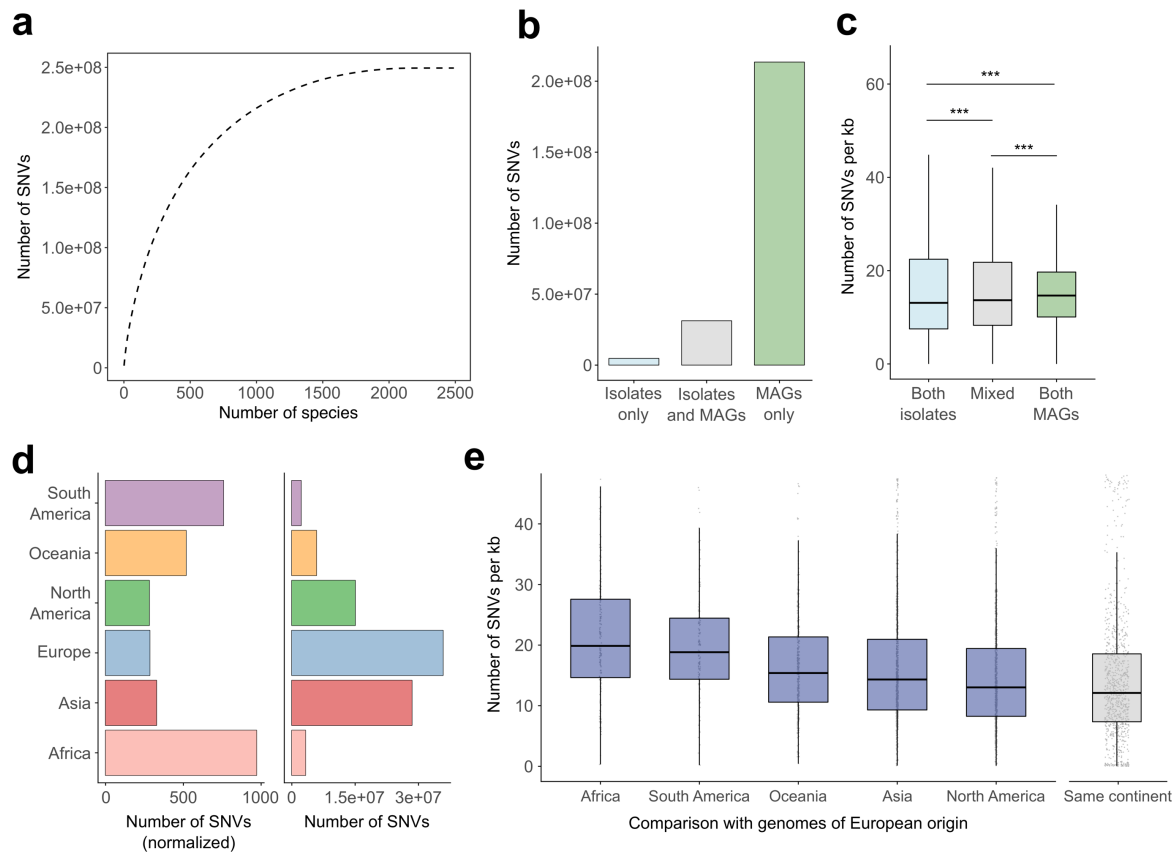
730 accessory genes from each species that was classified with various annotation schemes,

731 alongside the percentage of genes lacking any functional annotation. \*\*\* $P < 0.001$  **d**,

732 Comparison between the functional categories assigned to the core and accessory genes. Only

733 those statistically significant (adjusted  $P < 0.05$ ) are shown. A positive effect size indicates

734 overrepresentation in the core genes.

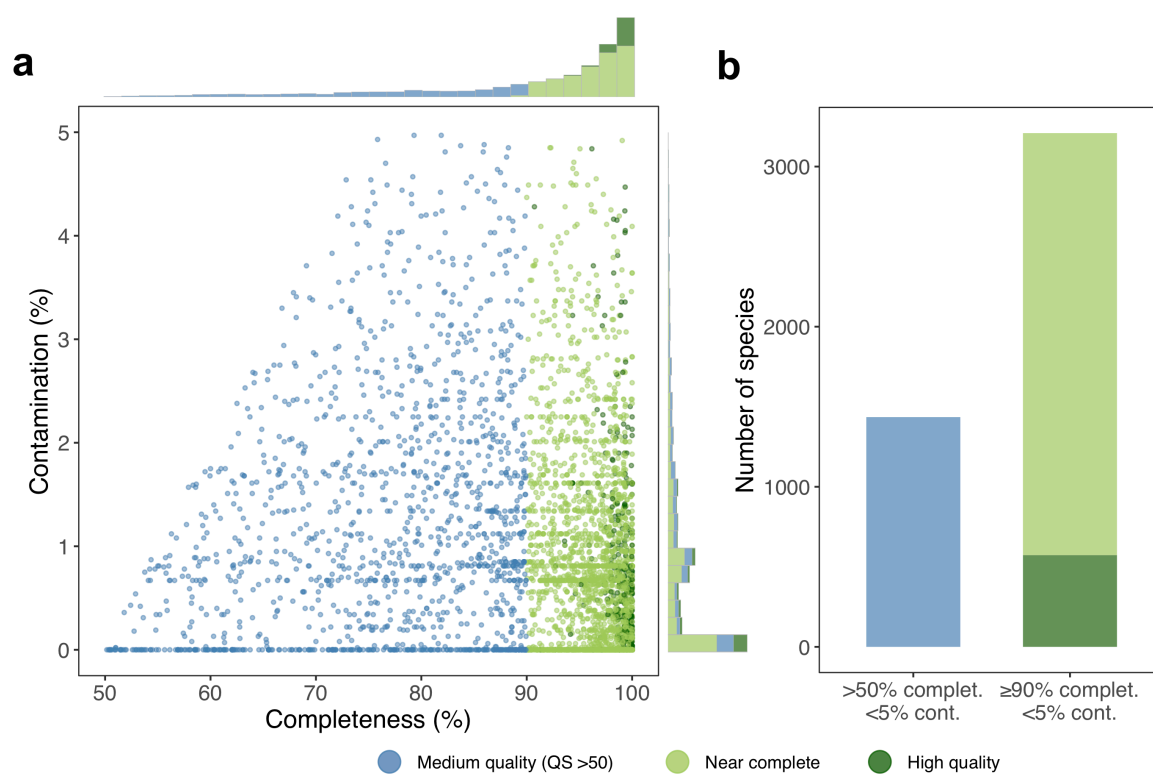


735

736 **Figure 6. Analysis of intra-species single nucleotide variation.**

737 **a**, Total number of SNVs detected as a function of the number of species. The cumulative  
 738 distribution was calculated after ordering the species by decreasing number of SNVs. **b**,  
 739 Number of SNVs detected only in isolate genomes, MAGs, or in both. **c**, Pairwise SNV density  
 740 analysis of genomes of the same or different type. \*\*\* $P < 0.001$  **d**, Right panel shows the  
 741 number of SNVs exclusively detected in genomes from each continent. The left panel shows  
 742 the number of exclusive SNVs normalized by the number of genomes per continent. **e**, Pairwise  
 743 SNV density analysis between genomes from Europe, the largest genome subset, and other  
 744 continents. The median SNV density was calculated per species and the distribution is shown  
 745 for all species. Comparison of genomes recovered from the same continent was used as  
 746 reference. The SNV density between genomes of the same continent is significantly lower  
 747 (adjusted  $P < 0.05$ ) to that calculated for genomes from different continents.

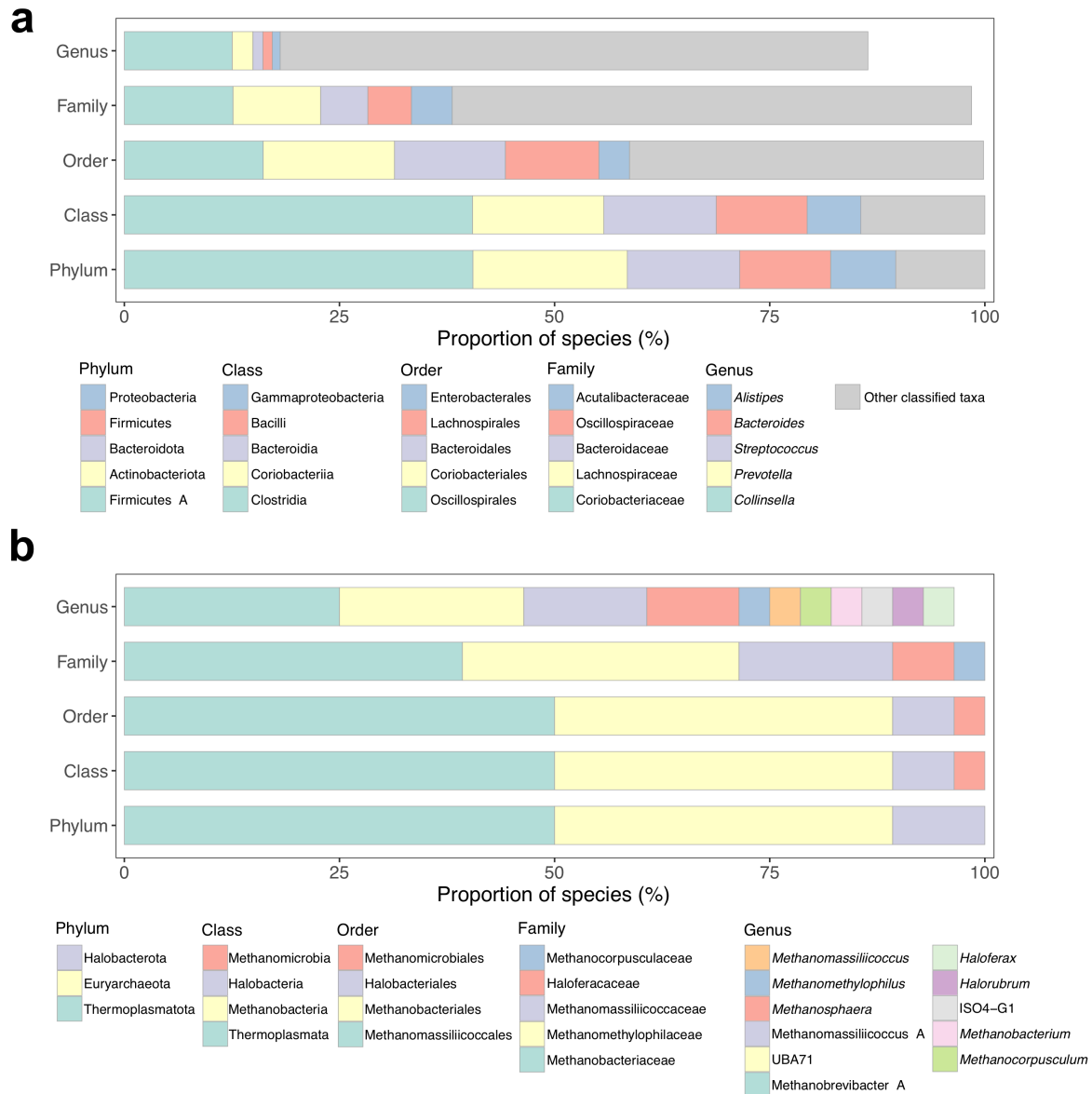
## 748 Supplementary Figures



749

### 750 Supplementary Figure 1. Genome quality of species representatives.

751 **a**, Completeness and contamination scores for each of the 4,644 species representatives,  
752 coloured by their quality classification category. Medium quality: >50% completeness; near  
753 complete: ≥90% completeness; high-quality: >90% completeness, presence of 5S, 16S and 23S  
754 rRNA genes, as well as at least 18 tRNAs. All genomes have a quality score (QS =  
755 completeness – 5 × contamination) above 50. **b**, Number of species according to different  
756 completeness and contamination criteria.



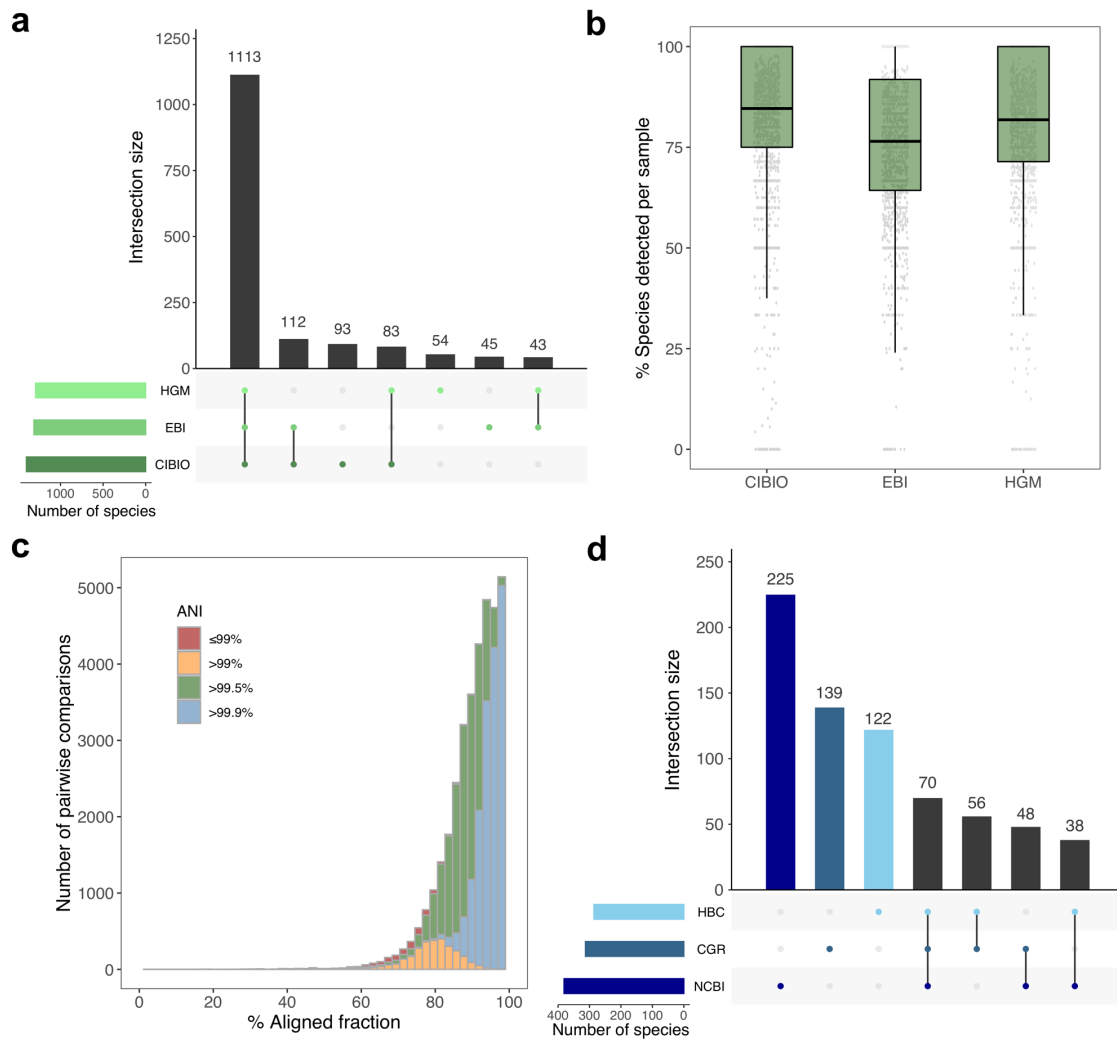
757

758 **Supplementary Figure 2. Taxonomy composition of the bacterial and archaeal species.**

759 **a**, Taxonomic affiliation of the 4,616 bacterial species detected. Data is partitioned by

760 taxonomic rank, with only the five most highly represented taxa per rank depicted in the legend.

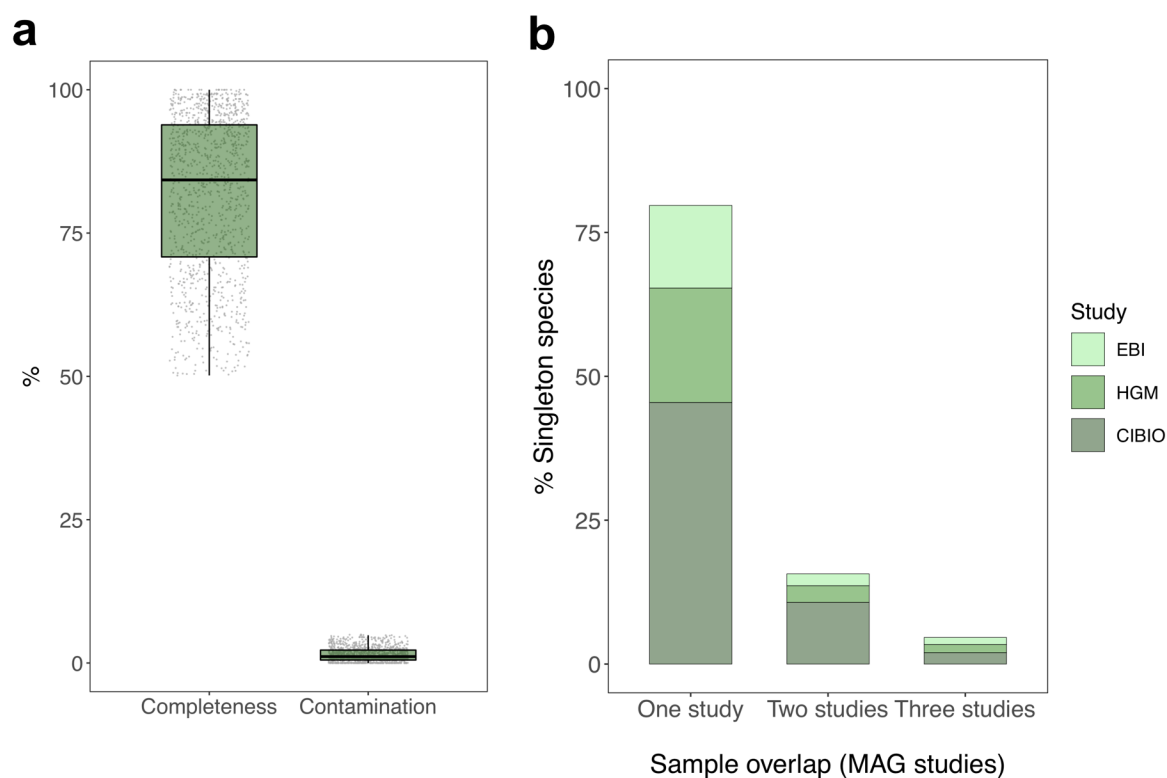
761 **b**, Taxonomic affiliation of the 28 archaeal species detected, partitioned by taxonomic rank.



762

763 **Supplementary Figure 3. Species overlap across study sets.**

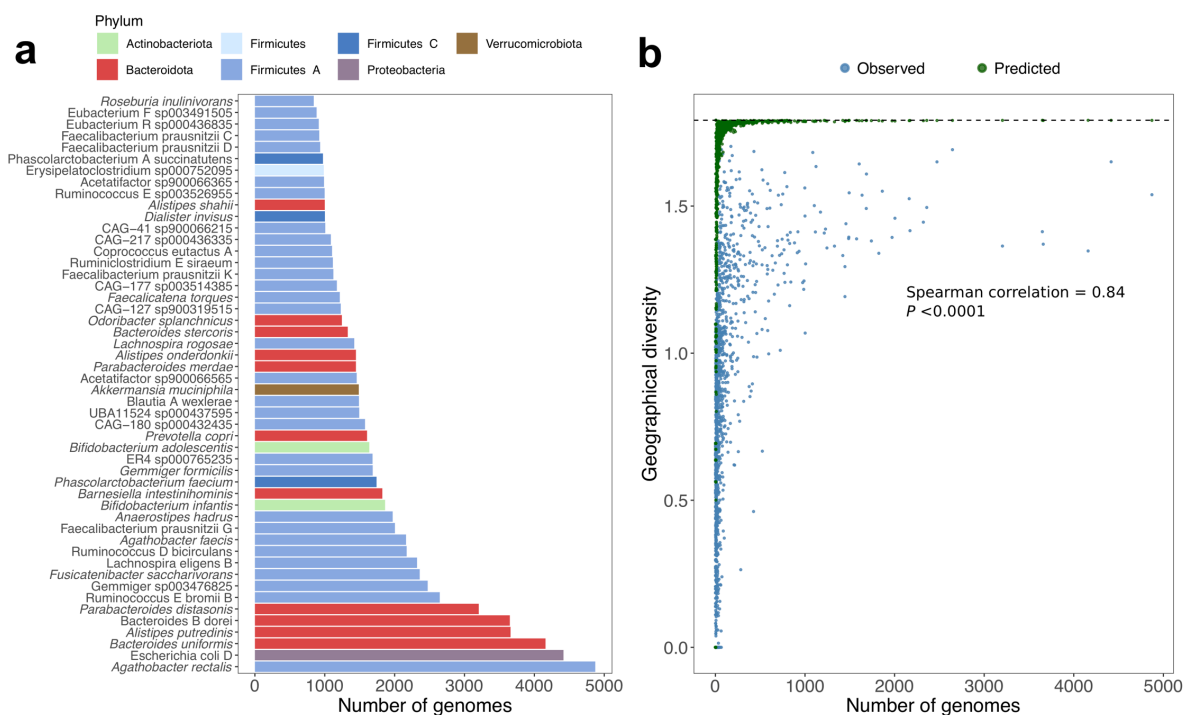
764 **a**, Number of species found across the three metagenome-assembled genome sets, ordered by  
 765 their level of overlap. Only those genomes recovered from the 1,554 metagenomic samples  
 766 used by all three studies were considered in this analysis. **b**, Distribution of the proportion of  
 767 species recovered per sample in each study out of all species recovered across all three studies  
 768 in the same samples. **c**, Estimated aligned fractions and average nucleotide identities (ANI)  
 769 between conspecific genomes obtained in the same sample but in different MAG studies. **d**,  
 770 Number of species identified in three culture-based studies and their degree of overlap. The  
 771 NCBI study set consists mainly of genomes from the Human Microbiome Project (HMP).



772

773 **Supplementary Figure 4. Quality and sample origin of uncultured singleton species.**

774 **a**, Genome completeness and contamination estimates of the 1,212 uncultured species  
775 represented by a single genome. **b**, Proportion of the 1,212 singleton species, by study set, that  
776 originated from samples analysed in one, two or three of the MAG studies (CIBIO, EBI and  
777 HGM). The CIBIO study used metaSPAdes and MetaBAT 2 for assembling and binning  
778 sequencing runs previously merged by sample; the HGM study used MEGAHIT to assemble  
779 runs merged by sample and applied a combination of MaxBin 2, MetaBAT 2, CONCOCT and  
780 DAS Tool for binning and refinement; the EBI study used metaSPAdes and MetaBAT 2 for  
781 assembling and binning individual runs without merging by sample.



782

783 **Supplementary Figure 5. Species frequency and geographical diversity.**

784 **a**, Number of non-redundant genomes retrieved from the 50 most highly represented species

785 in the UHGG. Each species is coloured by its assigned phylum according to the figure legend.

786 **b**, Geographical diversity estimated using the Shannon index in relation to the number of non-

787 redundant genomes from each species. The Spearman's rank correlation coefficient and *P* value

788 are depicted in the graph. Predicted values represent the random geographical distribution of

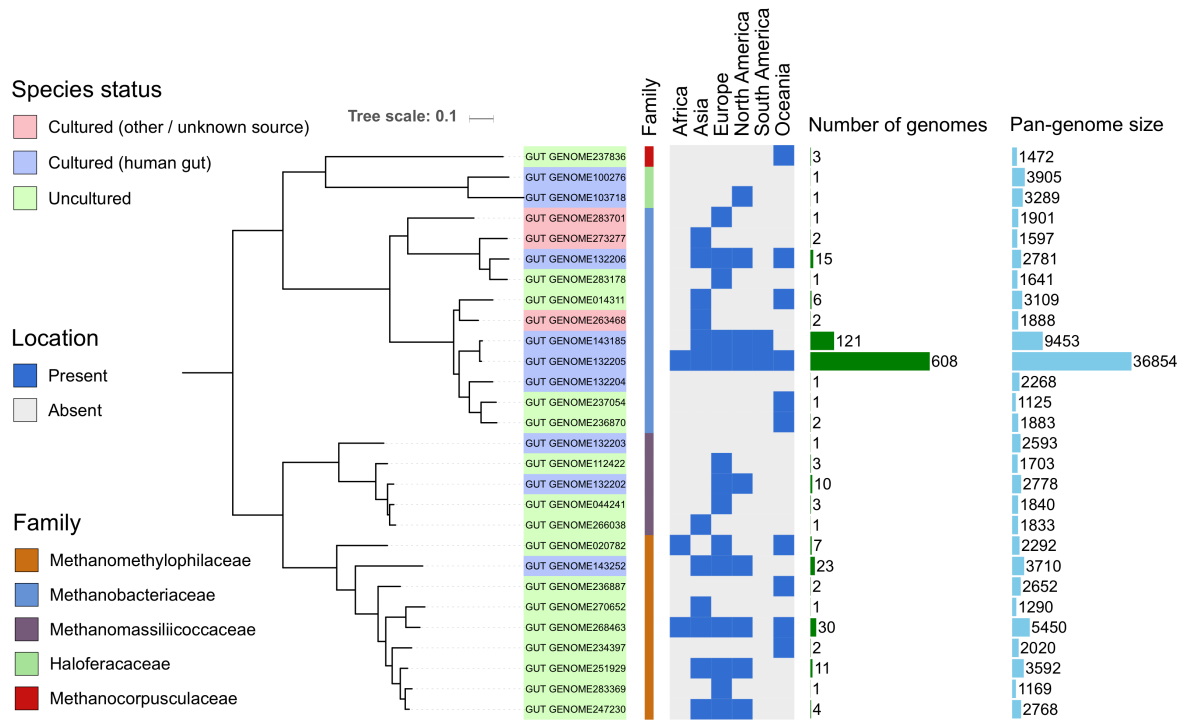
789 equivalent numbers of genomes observed for each species. Dashed horizontal line indicates the

790 maximum theoretical value of geographical diversity corresponding to equal sample

791 proportions across the six major continents (Africa, Asia, Europe, North America, South

792 America and Oceania).

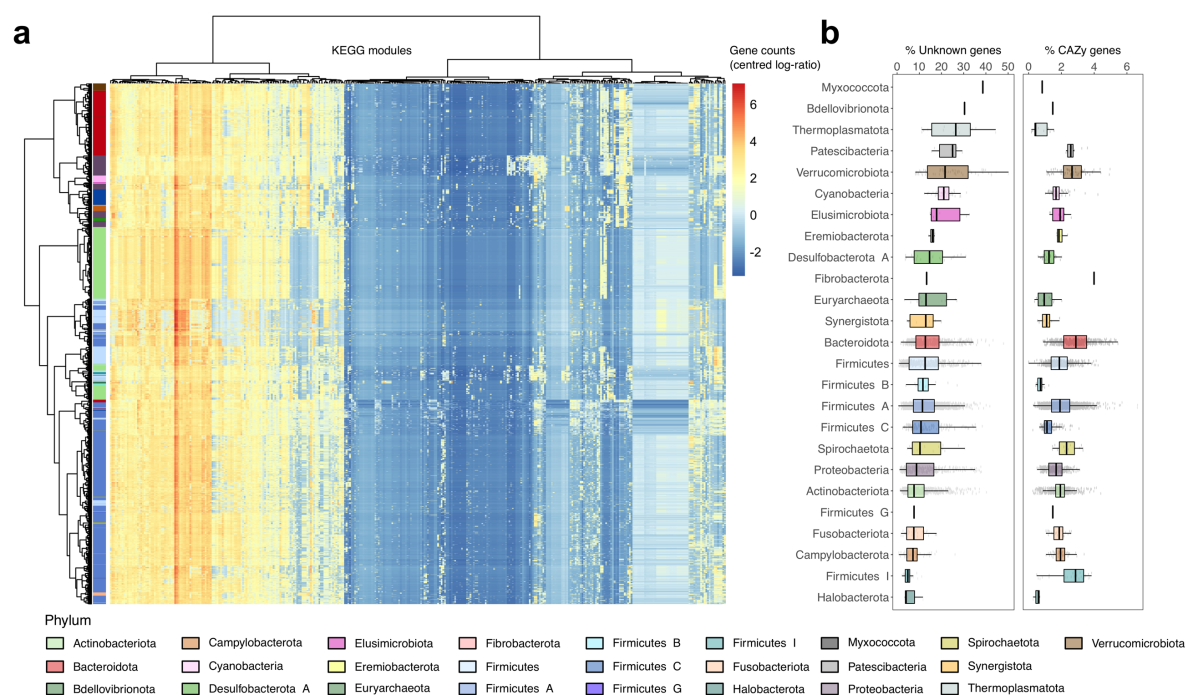




793

794 **Supplementary Figure 6. Diversity of the gut archaeal species detected.**

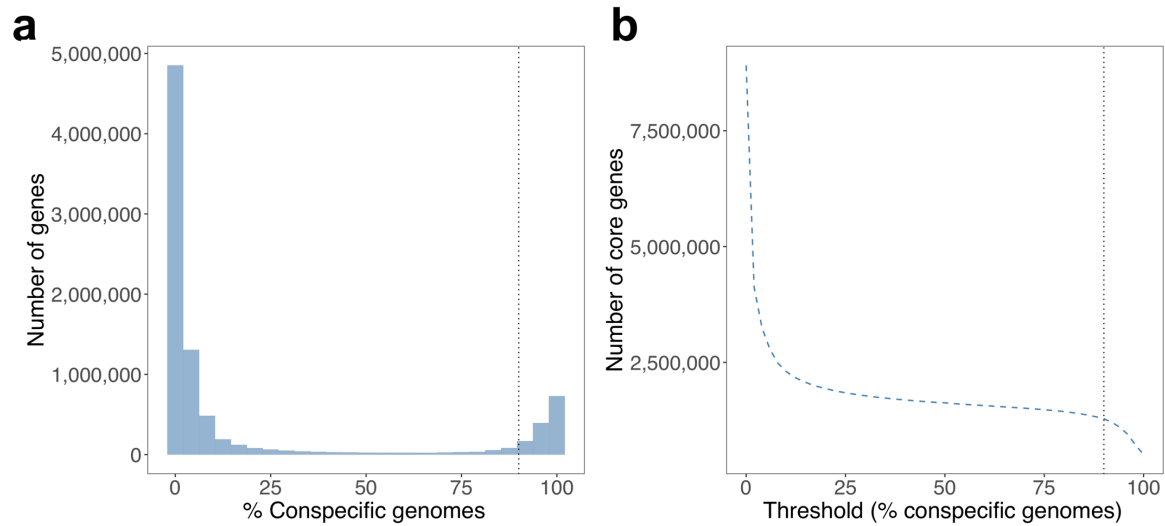
795 Phylogenetic tree of the 28 archaeal species detected in the human gut. Tips are labelled with  
 796 the corresponding species representative code and coloured according to its cultured status.  
 797 The taxonomic affiliation (family), geographical distribution, number of non-redundant  
 798 genomes and total pan-genome size are represented next to the tree.



799

800 **Supplementary Figure 7. Functional annotation of gut microbiome species.**

801 **a**, Functional profiles of the UHGG species pan-genomes (rows) according to 363 KEGG  
 802 modules (columns). Numbers of genes matching each module were normalized to centred log-  
 803 ratios after imputing values with zero counts. Species are coloured according to phylum.  
 804 KEGG modules and species were hierarchically clustered using the Ward's criterion method.  
 805 **b**, Proportion of each species pan-genome, partitioned by phylum, without any assignment to  
 806 the eggNOG, InterPro, COG or KEGG databases (left). Proportion of the pan-genome with a  
 807 match to the carbohydrate-active enzymes (CAZy) database (right).



808

809 **Supplementary Figure 8. Gene frequency distribution within the species-level clusters.**

810 **a**, Distribution of the number of genes found per fraction of conspecific genomes. Only near-  
811 complete genomes ( $\geq 90\%$  completeness) were considered in the analysis. **b**, Number of core  
812 genes detected based on the threshold of genomes per species used to classify as core. Vertical  
813 dashed line represents the 90% threshold used in this study.