

1 **Large scale metagenome assembly reveals novel animal-associated microbial genomes,**  
2 **biosynthetic gene clusters, and other genetic diversity**

3

4 Nicholas D. Youngblut<sup>\*,1</sup>, Jacobo de la Cuesta-Zuluaga<sup>1</sup>, Georg H. Reischer<sup>2,3</sup>, Silke  
5 Dauser<sup>1</sup>, Nathalie Schuster<sup>2</sup>, Chris Walzer<sup>4</sup>, Gabrielle Stalder<sup>4</sup>, Andreas H.  
6 Farnleitner<sup>2,3,5</sup>, Ruth E. Ley<sup>1</sup>

7 <sup>1</sup>Department of Microbiome Science, Max Planck Institute for Developmental Biology, Max Planck Ring 5,  
8 72076 Tübingen, Germany

9 <sup>2</sup>TU Wien, Institute of Chemical, Environmental and Bioscience Engineering, Research Group for  
10 Environmental Microbiology and Molecular Diagnostics 166/5/3, Gumpendorfer Straße 1a, A-1060  
11 Vienna, Austria

12 <sup>3</sup>ICC Interuniversity Cooperation Centre Water & Health, 1160 Vienna, Austria

13 <sup>4</sup>Research Institute of Wildlife Ecology, University of Veterinary Medicine, Vienna, Austria.

14 <sup>5</sup>Research Division Water Quality and Health, Karl Landsteiner University for Health Sciences, 3500  
15 Krems an der Donau, Austria

16

17

18

19

20

21

22 \* Corresponding author: Nicholas Youngblut (nyoungblut@tuebingen.mpg.de)

23

24 **Running title:** Metagenome assembly reveals vertebrate microbiome diversity

25 **Key words:** animal microbiome, gut, metagenome assembly, novel diversity

26

27

## 28 **Abstract**

29           Large-scale metagenome assemblies of human microbiomes have produced a  
30 vast catalogue of previously unseen microbial genomes; however, comparatively few  
31 microbial genomes derive from other vertebrates. Here, we generated 5596  
32 metagenome-assembled genomes from the gut metagenomes of 180 predominantly  
33 wild animal species representing 5 classes, in addition to 14 existing animal gut  
34 metagenome datasets. The MAGs comprised 1522 species-level genome bins (SGBs);  
35 most of which were novel at the species, genus, or family levels, and the majority were  
36 enriched in host versus environment metagenomes. Many traits distinguished SGBs  
37 enriched in host or environmental biomes, including the number of antimicrobial  
38 resistance genes. We identified 1986 diverse biosynthetic gene clusters; only 23  
39 clustered with any MIBiG database references. Gene-based assembly revealed  
40 tremendous gene diversity, much of it host- or environment-specific. Our MAG and gene  
41 datasets greatly expand the microbial genome repertoire and provide a broad view of  
42 microbial adaptations to the vertebrate gut.

## 43 **Importance**

44           Microbiome studies on a select few mammalian species (*e.g.*, humans, mice, and  
45 cattle) have revealed a great deal of novel genomic diversity in the gut microbiome.  
46 However, little is known of the microbial diversity in the gut of other vertebrates. We  
47 studied the gut microbiome of a large set of mostly wild animal species consisting of  
48 mammals, birds, reptiles, amphibians, and fish. Unfortunately, we found that existing  
49 reference databases commonly used for metagenomic analyses failed to capture the  
50 microbiome diversity among vertebrates. To increase database representation, we  
51 applied advanced metagenome assembly methods to our animal gut data and to many  
52 public gut metagenome datasets that had not been used to obtain microbial genomes.  
53 Our resulting genome and gene cluster collections comprised a great deal of novel  
54 taxonomic and genomic diversity, which we extensively characterized. Our findings  
55 substantially expand what is known of microbial genomic diversity in the vertebrate gut.

## 56 **Introduction**

57           The vertebrate gut microbiome comprises a vast amount of genetic diversity, yet  
58 even for the most well-studied species such as humans, the number of microbial  
59 species lacking a reference genome was recently estimated to be 40-50%<sup>1</sup>. Uncovering  
60 this “microbial dark matter” is essential to understanding the roles of individual  
61 microbes, their intra- and inter-species diversity within and across host populations, and  
62 how each microbe interacts with each other and the host to mediate host physiology in  
63 a myriad number of ways<sup>2</sup>. On a more applied level, characterizing novel gut microbial  
64 diversity aids in bioprospecting of novel bioactive natural products, catalytic and

65 carbohydrate-binding enzymes, probiotics, etc., along with aiding in the discovery and  
66 tracking of novel pathogens and antimicrobial resistance (AMR)<sup>3</sup>.

67 Recent advances in culturomic approaches have generated thousands of novel  
68 microbial genomes<sup>4-6</sup>, but the throughput is currently far outpaced by metagenome  
69 assembly approaches<sup>7</sup>. However, such large-scale metagenome assembly-based  
70 approaches have not been as extensively applied to most non-human vertebrates. The  
71 low amount of metagenome reads classified in some recent studies of the rhinoceros,  
72 chicken, cod, and cow gut/rumen microbiome suggests that databases lack much of the  
73 genomic diversity in less-studied vertebrates<sup>8-11</sup>. Indeed, the limited number of studies  
74 incorporating metagenome assembly hint at the extensive amounts of as-of-yet novel  
75 microbial diversity across the >66,000 vertebrate species on our planet.

76 Here, we developed an extensive metagenome assembly pipeline and applied it  
77 to a multi-species dataset of microbiome diversity across vertebrate species comprising  
78 5 classes: Mammalia, Aves, Reptilia, Amphibia, and Actinopterygii, with >80% of  
79 samples obtained from wild individuals<sup>12</sup> combined with data from 14 published animal  
80 gut metagenomes. Moreover, we also applied a recently developed gene-based  
81 metagenome assembly pipeline to the entire dataset in order to obtain gene-level  
82 diversity for rarer taxa that would otherwise be missed by genome-based assembly<sup>13,14</sup>.  
83 Our assembly approaches yielded a great deal of novel genetic diversity, which we  
84 found to be largely enriched in animals versus the environment, and to some degree,  
85 enriched in particular animal clades.

## 86 **Methods**

### 87 *Sample collection*

88 Sample collection was as described in Youngblut and colleagues<sup>12</sup>. Table S1A shows  
89 the dates, locations, and additional metadata of all samples collected. All fecal samples  
90 were collected in sterile sampling vials, transported to a laboratory and frozen within 8  
91 hours. DNA extraction was performed with the PowerSoil DNA Isolation Kit (MoBio  
92 Laboratories, Carlsbad, USA).

### 93 *“multi-species” vertebrate gut metagenomes*

94 Metagenome libraries were prepared as described by Karasov and colleagues<sup>15</sup>.  
95 Briefly, 1 ng of input gDNA was used for Nextera Tn5 tagmentation. A BluePippin was  
96 used to restrict fragment sizes to 400-700 bp. Barcoded samples were pooled and  
97 sequenced on an Illumina HiSeq3000 with 2x150 paired-end sequencing. Read quality  
98 control (QC) is described in the Supplemental Methods.

99 Post-QC reads were taxonomically profiled with Kraken2 and Bracken v.2.2<sup>16</sup>  
100 against the Struo-generated GTDB-r89 Kraken2 and Bracken databases<sup>17</sup>. HUMAnN2

101 v.0.11.2<sup>18</sup> was used to profile genes and pathways against the Struo-generated  
102 HUMAnN2 database created from GTDB-r89.

### 103 *Publicly available animal gut metagenomes*

104 Published animal gut metagenome reads were downloaded from the Sequence  
105 Read Archive (SRA) between May and August of 2019. Table S1B lists all included  
106 studies. We selected studies with Illumina paired-end metagenomes from gut contents  
107 or feces. MGnify samples were downloaded from the SRA in Oct 2019 (Table S1C).  
108 Read quality control is described in the Supplemental Methods.

### 109 *Metagenome assembly of genomes pipeline*

110 Assemblies were performed on a per-sample basis, with reads subsampled via  
111 seqtk v.1.3 to  $\leq 20$  million read pairs. The details of the assembly pipeline are described  
112 in the Supplemental Methods.

113 A multi-locus phylogeny of all SGB representatives was inferred with PhyloPhlAn  
114 v.0.41<sup>19</sup>. Secondary metabolites were identified with AntiSMASH v.5.1.1<sup>20</sup> and  
115 DeepBGC v.0.1.18<sup>21</sup> and then characterized with BiGSCAPE<sup>22</sup>. Abricate was used to  
116 identify antimicrobial resistance genes. We used Krakenuniq v.0.5.8<sup>23</sup> for estimating  
117 abundance of MAGs in metagenome samples. Details can be found in the  
118 Supplemental Methods.

### 119 *Metagenome assembly of genes pipeline*

120 Assemblies performed on a per-sample basis, with reads subsampled via seqtk  
121 v.1.3 to  $\leq 20$  million pairs. We used PLASS v.2.c7e35<sup>14</sup> and Linclust (mmseqs  
122 v.10.6d92c)<sup>13</sup> to assemble and cluster contigs. A full description is in the Supplemental  
123 Methods. DESeq2<sup>24</sup> was used to estimate enrichment of MAGs and gene clusters in  
124 metagenomes from host and environment biomes.

### 125 *Data availability*

126 The raw sequence data are available from the European Nucleotide Archive  
127 under the study accession number PRJEB38078. Fasta files for the 5596 non-  
128 redundant MAGs, 1522 SGBs, and gene clusters (50, 90, and 100% sequence identity  
129 clustering) can be found at  
130 [http://ftp.tue.mpg.de/ebio/projects/animal\\_gut\\_metagenome\\_assembly/](http://ftp.tue.mpg.de/ebio/projects/animal_gut_metagenome_assembly/), along with  
131 genbank files for all BGCs. Code used for processing the data can be found at  
132 [https://github.com/leylabmpi/animal\\_gut\\_metagenome\\_assembly](https://github.com/leylabmpi/animal_gut_metagenome_assembly).

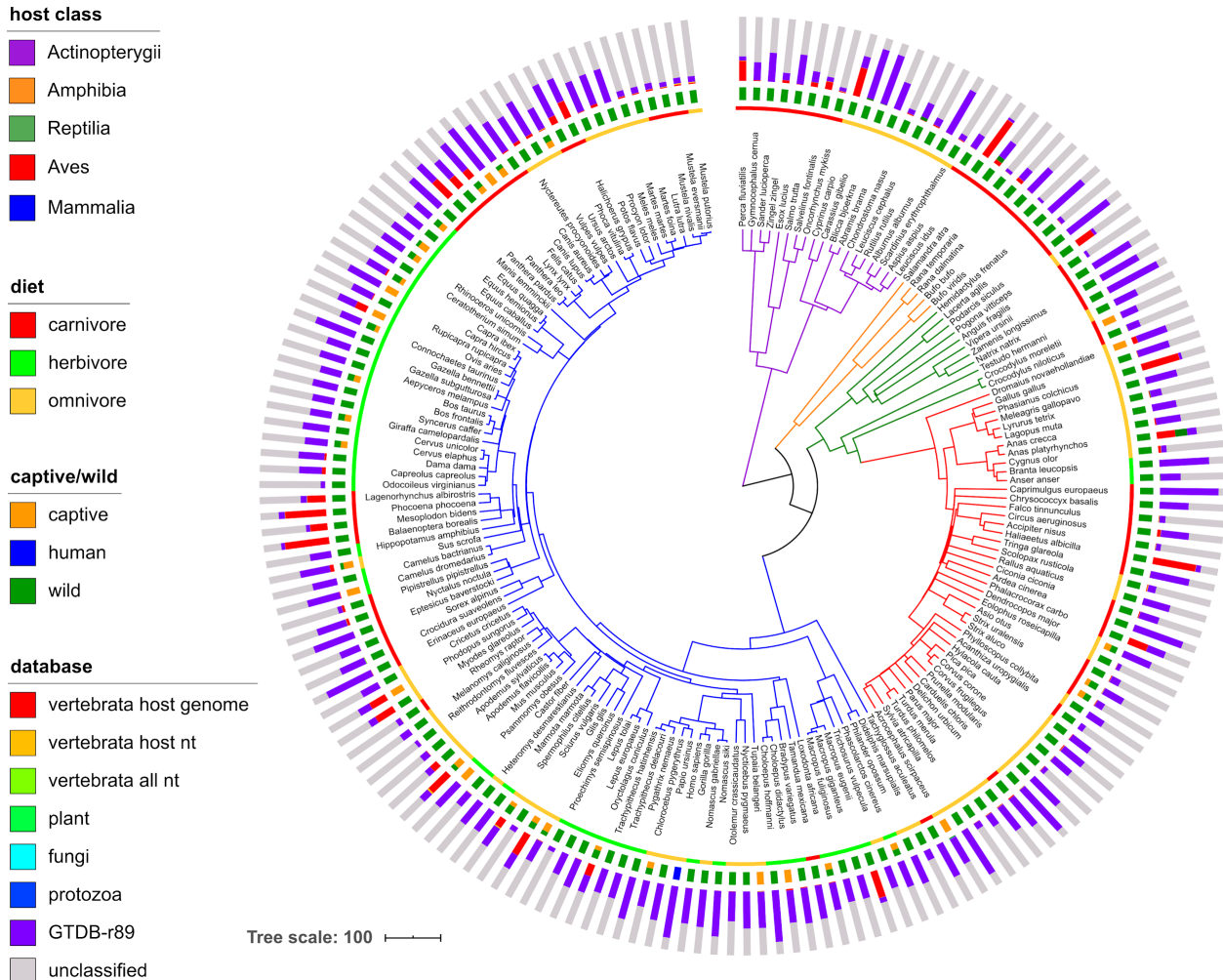
## 133 **Results**

### 134 *Animal gut metagenomes from a highly diverse collection of animals*

135 We generated animal gut metagenomes from a breadth of vertebrate diversity  
136 spanning five classes: Mammalia, Aves, Reptilia, Amphibia, and Actinopterygii (the  
137 “multi-species” dataset; Figure 1). In total, 289 samples passed our read quality control,  
138 with  $3.4e6 \pm 5e6$  s.d. paired-end reads per sample, resulting in a mean estimated  
139 coverage of  $0.54 \pm 0.14$  s.d. (Figure S1). 180 animal species were represented, with up  
140 to 6 individuals per species (mean of 1.6). Most individuals were wild (81%).

141 Our read-quality control pipeline included stringent filtering of host reads; some  
142 samples contained high amounts of reads mapping to vertebrate genomes (up to 74%;  
143  $6 \pm 17\%$  s.d.; Figure 1). Gut content samples contained a significantly higher amount of  
144 host reads ( $13.5 \pm 21.6\%$  s.d.) versus feces metagenomes ( $4.7 \pm 12.7\%$  s.d.; Wilcox,  $P$   
145  $< 1.8e-7$ ; Table S1A). We mapped all remaining reads to a custom comprehensive  
146 Kraken2 database built from the GTDB (Release 89). Still, many samples had a low  
147 percentage of mapped reads ( $43 \pm 22$  s.d.; Figure 1), with 29% of the samples having  
148  $<20\%$  mapped reads.

149



150  
 151 **Figure 1.** A large percentage of unmapped reads, even when using multiple comprehensive metagenome  
 152 profiling databases. The dated host species phylogeny was obtained from <http://timetree.org>, with  
 153 branches colored by host class. From inner to outer, the data mapped onto the tree is host diet, host  
 154 captive/wild status, and the mean number of metagenome reads mapped to various host-specific, non-  
 155 microbial, and microbial databases. Note that captive/wild status sometimes differs among individuals of  
 156 the same species. The databases are i) a representative of each publicly available genome from the host  
 157 species ("vertebrata host genome"), ii) all entries in the NCBI nt database with taxonomy IDs matching  
 158 host species ("vertebrata host nt"), iii) as the previous, but all vertebrata sequences included, iv) the  
 159 Kraken2 "plant" database, v) the Kraken2 "fungi" database, vi) the Kraken2 "protozoa" database, vii) a  
 160 custom bacteria and archaea database created from the Genome Taxonomy Database, Release 89  
 161 ("GTDB-r89"). Reads were mapped iteratively to each database in the order shown in the legend (top to  
 162 bottom), with only unmapped reads included in the next iteration. "unclassified" reads did not map to any  
 163 database, which were used along with reads mapping to GTDB-r89 for downstream analyses ("microbial  
 164 + unclassified").

165 *Discovery of novel diversity by large-scale metagenome assembly*

166 Our comprehensive metagenome assembly pipeline generated 4374 non-  
 167 redundant MAGs. After quality control and de-replication (see Methods), 296 MAGs  
 168 remained, with a mean percent completeness and contamination of  $84 \pm 14$  and  $1.5 \pm$   
 169  $1.2$  s.d., respectively (Figure S2; Supplemental Results).

170 We expanded our MAG dataset by applying our assembly pipeline to 14  
171 publically available animal gut metagenome datasets in which no MAGs have been  
172 generated by *de novo* metagenome assembly (Table S1B). Our metagenome selection  
173 included 554 samples from members of Mammalia (dogs, cats, woodrats, pigs, whales,  
174 rhinoceroses, pangolins, and non-human primates), Aves (geese, kakapos, and  
175 chickens), and Actinopterygii (cod). We applied our assembly pipeline to each individual  
176 dataset and generated a total of 5301 non-redundant MAGs (Figure S3; Supplemental  
177 Results). The substantially higher number of MAGs from these 14 datasets versus our  
178 single multi-species dataset is likely due to the larger number of samples and the high  
179 sequencing depth for many of those samples (*e.g.*, we used 2 billion paired-end reads  
180 from the dog gut microbiome dataset<sup>25</sup>).

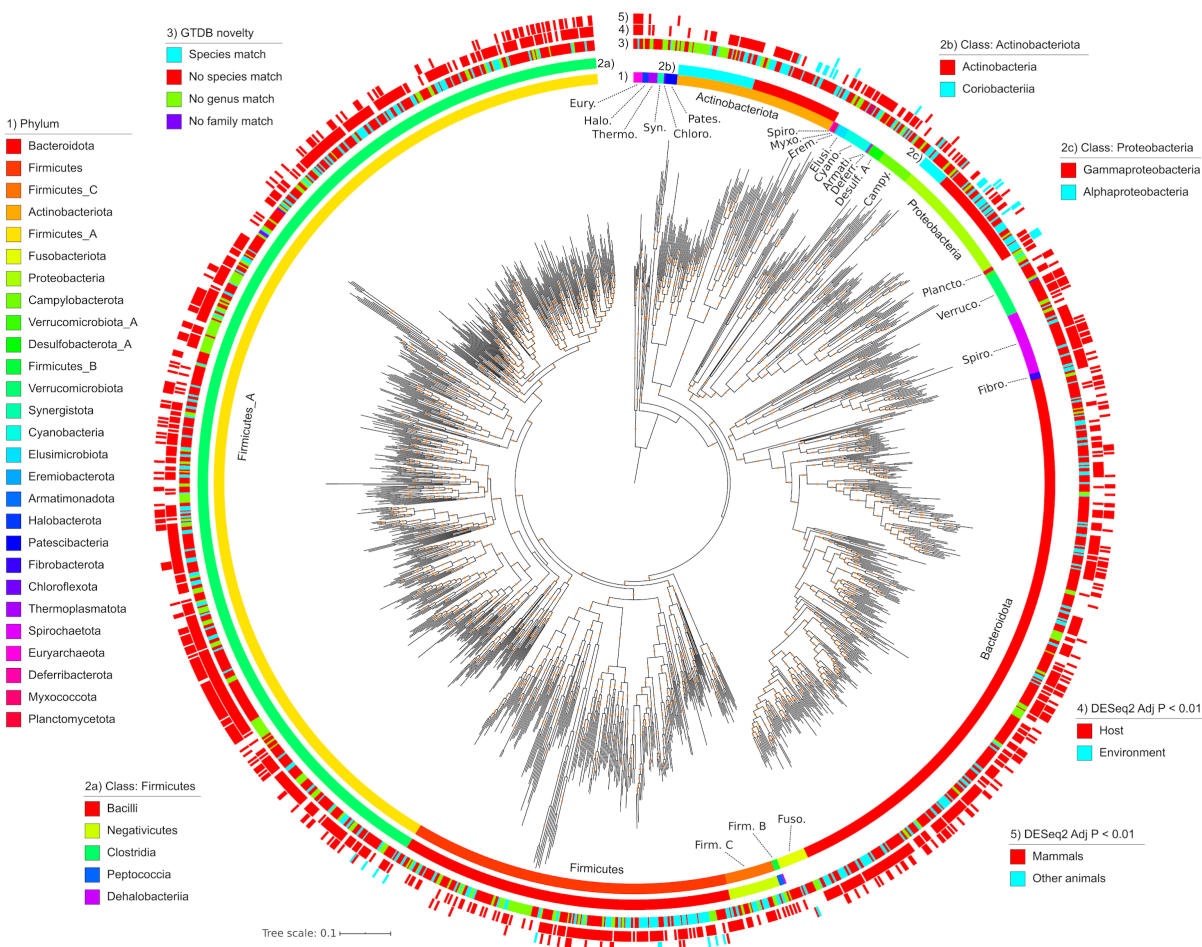
181 We combined all MAGs and de-replicated at 99.9 and 95% average nucleotide  
182 identity (ANI) to produce 5596 non-redundant MAGs and 1522 species-level genome  
183 bins (SGBs), respectively (Tables S2A & S2B). Of the 5596 MAGs, 2773 (50%) had a  
184 completeness of  $\geq 90\%$ . Of the 1522 SGBs, 1184 (78%) lacked a  $\geq 95\%$  ANI match to  
185 the GTDB-r89, 266 (17%) lacked a genus-level match, and 6 lacked a family-level  
186 match (Figures 2 & S4). Mapping taxonomic novelty onto a multi-locus phylogeny of all  
187 1522 SGBs revealed that novel taxa were rather dispersed across the phylogeny  
188 (Figure 2).

189 We also assessed the novelty of our SGBs relative to UHGG, a comprehensive  
190 human gut genome database, and found that only 31% of our SGBs had  $\geq 95\%$  ANI to  
191 any of the 4644 UHGG representatives, and this overlap only increased to 34% at a  
192 90% ANI cutoff.

193 Our SGB collection mostly consisted of MAGs assembled from a few species in  
194 the multi-study dataset, suggesting that the SGBs may not be representative of taxa  
195 found in other, more distantly related vertebrates. To assess the level of representation,  
196 we determined the prevalence of all SGBs across all multi-species metagenomes  
197 (Figure S5). The host species with the highest number of observed SGBs did tend to be  
198 those comprising the multi-study dataset (*e.g.*, pigs and primates); however, SGBs were  
199 frequently observed across the host phylogeny ( $41 \pm 61$  s.d. SGBs per host), indicating  
200 that the SGB collection was generally representative of the vertebrate gut microbiome.

201 Integrating the 1522 SGBs into our custom GTDB Kraken2 database significantly  
202 increased the percent reads mapped (paired t-test,  $P < 0.005$ ; Figure S6). The percent  
203 increase varied from  $<1$  to 62.8% (mean of  $5.3 \pm 6.7$  s.d.) among animal species but did  
204 not appear biased to just pigs, dogs, or other vertebrate species in the multi-study  
205 datasets that we incorporated (Figure S7), which corresponds with our analysis of SGB  
206 prevalence across vertebrate hosts (Figure S5).

207



208  
 209 **Figure 2.** A phylogeny of all 1522 SGBs. From innermost to outermost, the data mapped onto the  
 210 phylogeny is: GTDB phylum-level taxonomic classifications, class-level taxonomies for Actinobacteriota,  
 211 class-level taxonomies for Firmicutes, class-level taxonomies for Proteobacteria, taxonomic novelty,  
 212 significant enrichment in host gut or environmental metagenomes, and significant enrichment in Mammals  
 213 or other animals in our multi-species gut metagenome dataset. The phylogeny was inferred from multiple  
 214 conserved loci via PhyloPhlAn. Orange dots on the phylogeny denote bootstrap values in the range of 0.7  
 215 to 1. The phylogeny is rooted on the last common ancestor of Archaea and Bacteria.

## 216 *Enrichment of SGBs among animal clades*

217 While the MAGs generated here derive from animal gut metagenomes, many of  
 218 these taxa might be transient in the host and actually more prevalent in the  
 219 environment. We tested this by generating a “host-environment” metagenome dataset  
 220 comprising 283 samples from 30 BioProjects (17 environmental and 13 host-associated;  
 221 Figure 3A). We found 932 of the 1522 SGBs (61%) to be significantly enriched in the  
 222 host metagenomes (DESeq2, *adj. P* < 0.01; Figure 3B). The host-enriched SGBs (host-  
 223 SGBs) were taxonomically diverse, comprising 22 phyla. In contrast, only 15 SGBs (1%)  
 224 were environment-enriched (env-SGBs), which all belonged to either Actinobacteriota or  
 225 Proteobacteria (Figure 3B). The only SGBs that were not significantly enriched in either  
 226 group belonged to Actinobacteriota or Proteobacteria, along with two SGBs from the



227 Firmicutes A phylum. Mapping these data onto the SGB phylogeny revealed  
228 phylogenetic clustering of the environment-enriched SGBs (Figure 2).

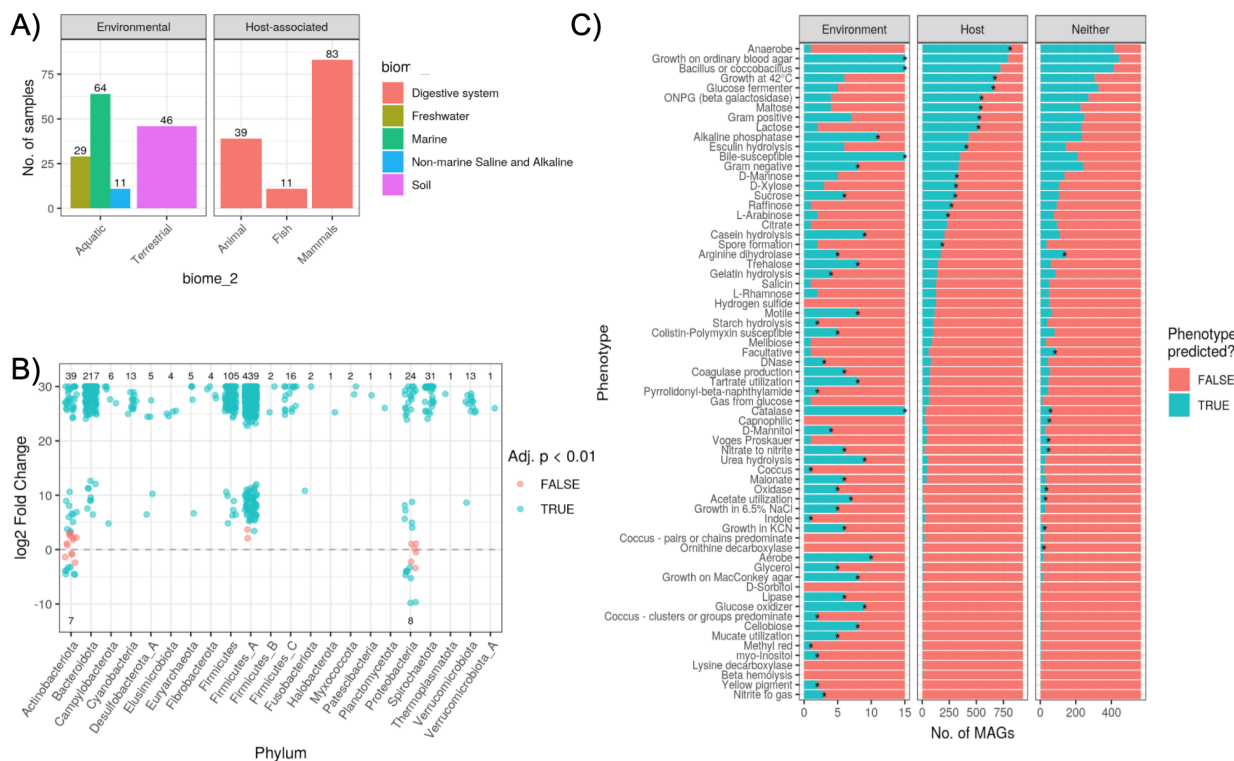
229 We investigated the traits of the host- and environment-enriched SGBs and  
230 found many predicted phenotypes to be more prevalent in one or the other group  
231 (Figure 3C; Table S2C). A total of 67 traits were predicted based on genomic content of  
232 certain pfam domains<sup>26</sup>. Almost all env-SGBs were aerobes (93%), which may aid in  
233 transmission between the environment and host biomes. In contrast, 87% of host-SGBs  
234 were anaerobes. Furthermore, all env-SGBs could generate catalase and were bile  
235 susceptible, while both phenotypes were sparse in host-SGBs (Figure 3C).  
236 Carbohydrate metabolism also differed, with most host-SGBs predicted to consume  
237 various tri-, di-, and mono-saccharides. In contrast, env-SGBs were enriched in  
238 phenotypes associated with motility, nitrogen metabolism, and breakdown of  
239 heterogeneous substrates (*e.g.*, cellobiose metabolism).

240 We also compared SGB enrichment in mammals versus non-mammals in our  
241 “multi-species” metagenome dataset and found 361 SGBs (24%) to be significantly  
242 enriched in mammals, while 22 (1%) were enriched in non-mammals (DESeq2, *adj. P* <  
243 0.01; Figure S2C. Interestingly, 100% of SGBs in the two archaeal phyla (Halobacteria  
244 and Euryarchaeota) were enriched in mammals. Also of note, most of the  
245 Verrucomicrobiota SGBs (87%) were enriched in mammals. The only 2 phyla with >10%  
246 of SGBs enriched in non-mammals were Proteobacteria (29%) and Campylobacteria  
247 (25%).

248 In contrast to our assessment of phenotypes distinct to host- or env-SGBs, we  
249 did not observe such a distinction of phenotypes among SGBs enriched in Mammalia or  
250 non-mammal gut metagenomes (Figure S8). Certain phenotypes such as anaerobic  
251 growth and lactose consumption were more prevalent among mammal species, but they  
252 were not found to be significantly enriched relative to the null model.

253 Little is known about the distribution of antimicrobial resistance genes in the gut  
254 microbiomes of most vertebrate species<sup>27</sup>; therefore, we investigated the distribution of  
255 AMR genes among MAGs enriched in the environment versus host biomes. We found a  
256 mean of  $35 \pm 26$  s.d. AMR markers per genome (Figure S9A). The high average was  
257 largely driven by Proteobacteria and Campylobacter genomes, which had a mean of  
258 387 and 161 AMR markers per genome, respectively. The 5 most abundant markers  
259 were *ruvB*, *galE*, *tupC*, *fabL* (*ygaA*), and *arsT* (Figure S9A). The more abundant  
260 markers predominantly belonged to Firmicutes A, while Proteobacteria comprised larger  
261 fractions of the less abundant markers. Environment-enriched taxa contained  
262 substantially more AMR genes than host-enriched taxa, and the same was true for non-  
263 Mammalia versus Mammalia-enriched taxa (Figures S9B & S9C).

264



265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275

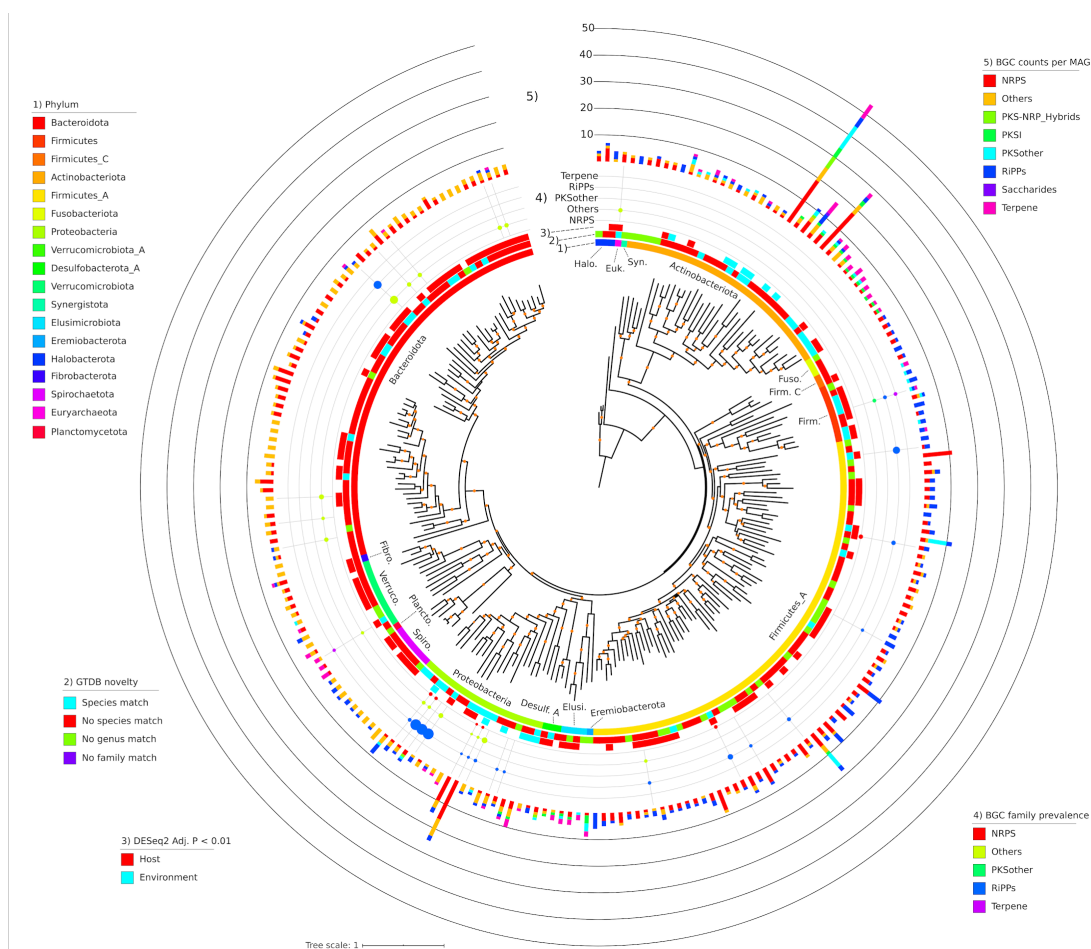
**Figure 3.** A) Summary of the number of samples per biome for our multi-environment metagenome dataset selected from the MGnify database. B) Number of SGBs found to be significantly enriched in relative abundances in host versus (positive log<sub>2</sub> fold change; "l2fc") environmental metagenomes (negative l2fc). Values shown are the number of MAGs significantly enriched (blue) in either biome or not found to be significant (red). C) Host- and environment-enriched SGBs have distinct traits. Phenotypes predicted based on MAG gene content (via TraitR<sup>26</sup>) are summarized for the SGBs significantly enriched in host or environmental metagenomes (DESeq2 Adj. P < 0.01) or neither biome ("Neither" in the x-axis facet). Note the difference in x-axis scale. Asterisks denote phenotypes significantly more prevalent in SGBs of the particular biome versus a null model of 1000 permutations in which biome labels were shuffled among SGBs. See Table S3A for all DESeq2 results.

### 276 MAGs reveal novel secondary metabolite diversity

277 We identified 1986 biosynthetic gene clusters (BGCs) among all 1522 SGBs. A  
278 total of 28 different products were predicted, with the most abundant being non-  
279 ribosomal peptide synthetases (NPRS; *n* = 473), sactipeptides (*n* = 307), and  
280 arylpolyenes (*n* = 291; Figure S10). BGCs were identified in 2 archaeal and 18 bacterial  
281 phyla. MAGs in the Firmicutes A phylum contained the most BGCs (*n* = 764; 38%),  
282 while Bacteroidota and Actinobacteriota phyla possessed 381 (19%) and 272 (14%),  
283 respectively (Figure S10). Still, Actinobacteriota SGBs did possess the highest average  
284 number of BGCs per genome (16.3), followed by Eremiobacterota (9), Proteobacteria  
285 (7.7), and Halobacterota (5.1).

286 Clustering all 1986 BGCs by BiGSCAPE generated 1764 families and 1305  
287 clans, with clans being a second, coarser level of clustering<sup>22</sup>. Only 8 clans (comprising  
288 23 BGCs) included any MIBiG database reference, suggesting a high degree of novelty

289 (Figure S11). Mapping the BGCs on a genome phylogeny of all species containing  $\geq 3$   
 290 BGCs (233 SGBs) revealed that the number of BGCs per genome was somewhat  
 291 phylogenetically clustered: the five genomes with the most BGCs belonged either to the  
 292 Actinobacteria or Gammaproteobacteria (Figure 4). Notably, these clades contained a  
 293 high number of host-SGBs. Of these 233 SGBs, the majority were taxonomically novel,  
 294 with 62% lacking a species-level match to GTDB-r89, and 18% lacking a genus-level  
 295 match (Figure 4). To determine which of the BGCs are most prevalent across animal  
 296 hosts, we quantified the prevalence of each BGC family across our multi-species  
 297 metagenome dataset and mapped it to the genome phylogeny (Figures 4 & S12). Of the  
 298 1543 BGC families found in the 233 SGBs, 83 were present in  $\geq 25\%$  of the animal  
 299 metagenomes, with ribosomally synthesized and post-translationally modified peptides  
 300 (RIPPs) being by far the most prevalent (up to 98% prevalence of individual BGC  
 301 families) and also found in species from a number of phyla.  
 302



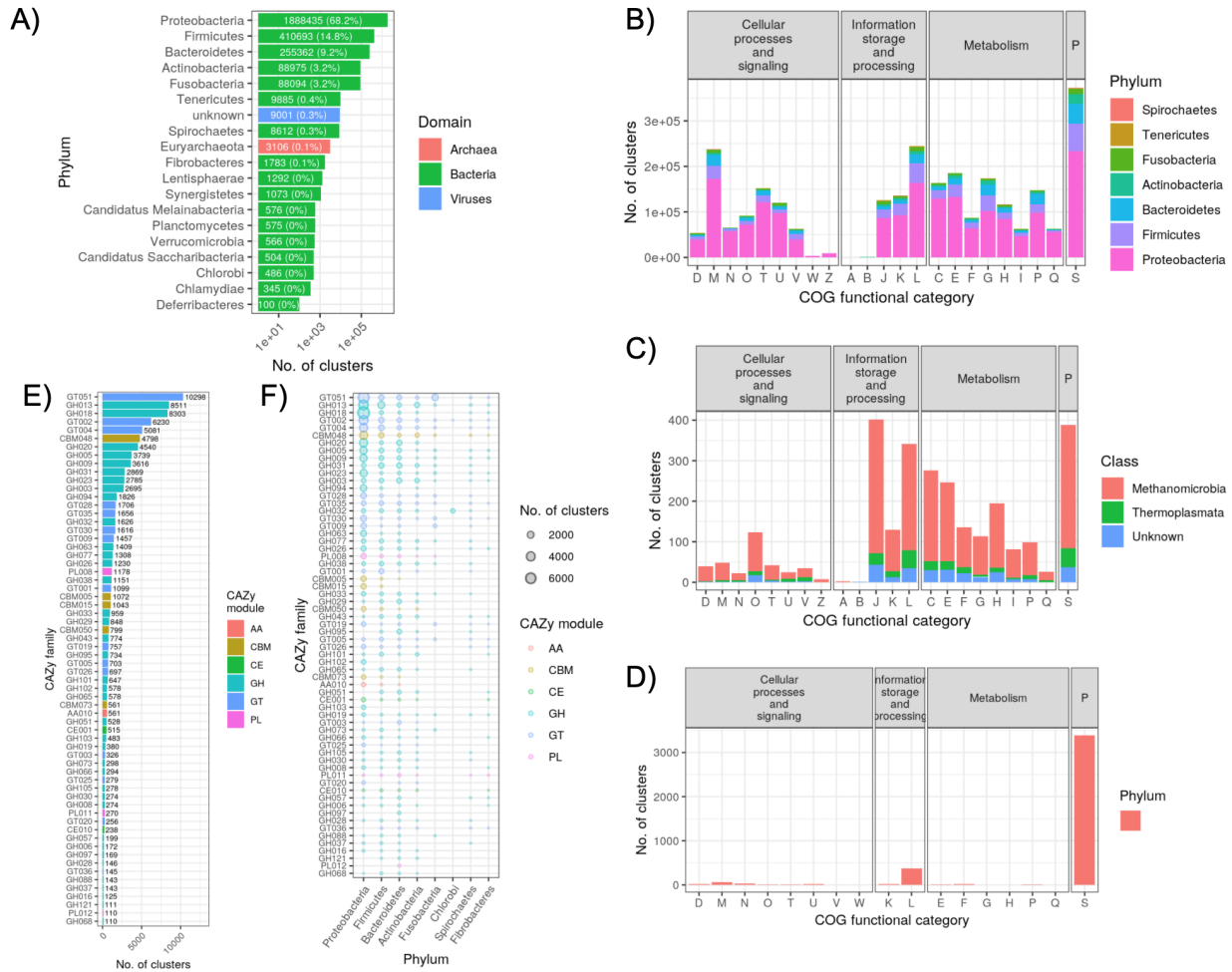
303  
 304 **Figure 4.** Phylogeny of all SGBs ( $n = 233$ ) with  $\geq 3$  BGCs identified by AntiSMASH. From innermost to  
 305 outermost, the data mapped onto the phylogeny is: 1) GTDB phylum-level taxonomic classifications, 2)  
 306 taxonomic novelty, 3) significant enrichment in host or environmental metagenomes, 4) the prevalence of  
 307 BGC families across the multi-species metagenome dataset, and 5) the number of BGCs identified in the  
 308 MAG. Prevalence is the maximum of any BGC family for that BGC type, and only BGC families with a

309 prevalence of  $\geq 25\%$  are shown. The phylogeny is a pruned version of that shown in Figure 2. Orange  
310 dots on the phylogeny denote bootstrap values in the range of 0.7 to 1.

### 311 *Large-scale gene-based metagenome assembly reveals novel diversity*

312 We applied gene-based assembly methods to our combined metagenome  
313 dataset<sup>14</sup>, which generated a total of 150,718,125 non-redundant coding sequences  
314 (average length of 179 amino acids). Clustering at 90 and 50% sequence identity  
315 resulted in 140,225,322 and 6,391,861 clusters, respectively. Only 16.9 and 11.3% of  
316 each respective cluster set mapped to the UniRef50 database, indicating that most  
317 coding sequences were novel. The clusters comprised 88 bacterial and 11 archaeal  
318 phyla; 80 of which were represented by  $< 100$  clusters, and 60 lacking a cultured  
319 representative. Proteobacteria (mostly Gammaproteobacteria), Firmicutes, and  
320 Bacteroidetes made up 92.2% of all clusters (Figure 5A). The proportion of clusters  
321 belonging to each COG functional category was largely the same for the more abundant  
322 bacterial phyla (Figure 5B), while more variation was seen among Euryarchaeota  
323 (Figure 5C). The dominant 7 phyla showed substantial variation in the number of  
324 clusters associated with various KEGG pathway categories (Figure S14). For instance,  
325 a high proportion of Fusobacteria and Tenericutes clusters were associated with the  
326 “nucleotide metabolism”, “replication and repair”, and “translation” categories. A total of  
327 87,573 clusters were annotated as CAZy families, with GT51, GH13, GH18, GT02, and  
328 GT04 representing 48% of all CAZy-annotated clusters (Figure 5E). Of the 12 phyla with  
329 the most CAZy family clusters, there were substantial differences in proportions of  
330 clusters falling into each family (Figure 5F).

331



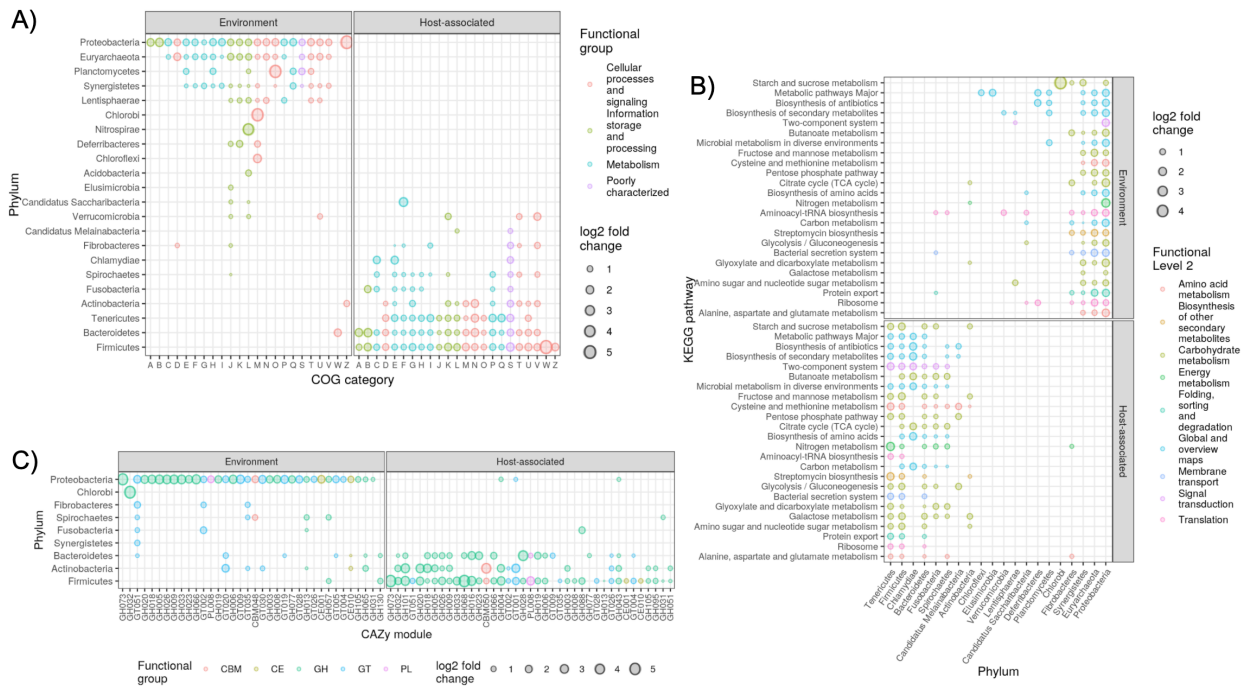
332  
 333 **Figure 5.** A summary of the 50% sequence identity clusters generated from the gene-based metagenome  
 334 assembly of the combined dataset. A) The total number of gene clusters per phylum. For clarity, only  
 335 phyla with  $\geq 100$  clusters are shown. Labels on each bar list the number of clusters (and percent of the  
 336 total). B) The number of bacterial gene clusters per phylum and COG category. The “P” facet label refers  
 337 to “poorly characterized”. C) The number of archaeal gene clusters per class (all belonging to  
 338 Euryarchaeota) and COG category. D) The number of viral gene clusters per COG category. E) The  
 339 number clusters annotated as each CAZy family. For clarity, only phyla with  $\geq 100$  clusters are shown.  
 340 Labels next to each bar denote the number of clusters. F) The number of clusters per CAZy family,  
 341 broken down by phylum. CAZy families and phyla are ordered by most to least number of clusters. For  
 342 clarity, only CAZy families and phyla with  $\geq 100$  total clusters are shown.

### 343 Biome enrichment of gene clusters from specific phyla

344 We mapped reads from our host-environment metagenome dataset to each  
 345 cluster and used DESeq2 to identify those significantly enriched ( $adj. P < 1e-5$ ) in each  
 346 biome. Most strikingly, the same functional groups were enriched in both biomes,  
 347 regardless of the grouping (*i.e.*, COG functional category, KEGG pathway, or CAZy  
 348 family); however, the gene clusters belonged to different microbial phyla (Figure 6;  
 349 Supplemental Results). For instance, nearly all COG categories for gene clusters  
 350 belonging to Proteobacteria were environment-enriched, while the same COG

351 categories for clusters belonging to Firmicutes and Bacteroidetes were host-enriched. In  
 352 contrast, functional groups of certain phyla were enriched in one biome, while different  
 353 groups were enriched in the other, indicating within-phylum differences in functional  
 354 content and habitat distributions. For instance, Fusobacteria KEGG pathways were  
 355 predominantly host-enriched, but protein export, bacteria secretion system, and  
 356 aminoacyl-tRNA biosynthesis were environment-enriched, indicating that these 3  
 357 pathways were more predominant in environment-enriched members of Fusobacteria  
 358 (Figure 6B). Overall, these results suggest that both biomes select for these same  
 359 microbial functions, but the microbes involved often differ at coarse taxonomic scales.

360 We also assessed gene cluster enrichment in Mammalia versus non-Mammalia  
 361 and found fewer significantly enriched features, which may be due to the smaller  
 362 metagenome sample size or less pronounced partitioning of functional groups among  
 363 biomes (Figure S15; Supplemental Results). Still, we again observed that both biomes  
 364 enriched for the same microbial functions, but these belonged to different coarse  
 365 taxonomic groups. To assess whether abundance estimations were substantially  
 366 erroneous due to mis-mapping of metagenome reads to the gene clusters, we reran the  
 367 analysis with stricter DIAMOND mapping parameters but observed similar findings,  
 368 even though 48% fewer gene clusters were detected in any metagenome (Figure S16).  
 369



370  
 371 **Figure 6.** Enrichment of gene clusters grouped by phylum and A) COG category B) KEGG pathway or C)  
 372 CAZy family. Only groupings significantly enriched in abundance (DESeq2, *adj. P* < 1e-5) in either biome  
 373 are shown. Only gene clusters observed in at least 25% of the metagenomes were included. For clarity,  
 374 only KEGG pathways enriched in >7 phyla are shown, and only CAZy families enriched in >1 phylum are  
 375 shown. Note that the axes are flipped in B) relative to A) and C). See Tables S5A, S5B, and S5C for all  
 376 DESeq2 results.

### 377 *Functional metagenome profiling benefits from our gene catalogue*

378         Lastly, we integrated our gene catalogue into a custom HUMAnN2 database built  
379 from the GTDB-r89 and found that this combined database substantially increased  
380 mappability of reads from our multi-species metagenome dataset (Figure S17;  
381 Supplemental Results).

## 382 **Discussion**

383         Our MAG and gene cluster datasets, derived from 289 newly generated  
384 metagenomes from 180 vertebrate species, along with 544 metagenomes from 14  
385 publicly available animal gut metagenome datasets, substantially helps to expand the  
386 breadth of cross-species gut metagenome comparisons (Figures 1 & 5). While  
387 metagenomics is rapidly expanding in popularity<sup>7</sup>, most analyses of metagenomic data  
388 suffer from a reliance on incomplete reference databases<sup>17</sup>, which we show to be  
389 acutely problematic for the gut microbiomes of most vertebrates in our dataset (Figure  
390 1). Grossly incomplete surveys of microbial diversity can lead to incorrect findings on  
391 community assembly in the vertebrate gut<sup>28</sup>. Although our dataset has only partially  
392 revealed this unknown diversity, it does substantially improve reference database  
393 coverage at both the genome and gene levels and also provides an estimate of the  
394 incompleteness of existing reference databases.

395         A major contribution of this study is the extensive MAG collection that we  
396 generated by assembling the metagenomes of our multi-species dataset together with  
397 14 other animal gut metagenome datasets from understudied host species. This  
398 collection includes 1184, 266, and 6 genomes from novel species, genera, and families,  
399 respectively (Figures 2 & S4). Moreover, we found little overlap (31%) between our  
400 MAG collection and the extensive human microbiome genome catalogue comprising the  
401 UHGG, which underscores its taxonomic novelty. We also showed substantial SGB  
402 prevalence across all 5 vertebrate taxonomic classes (Figure S5), indicating that our  
403 MAG collection is representative microbes found across the vertebrate taxonomy. Our  
404 MAG collection, once combined with the GTDB<sup>29</sup>, improved our ability to classify reads  
405 in our multi-species metagenome dataset (Figure S6), which is critical for accurately  
406 assessing gut microbiome diversity across vertebrates. Although MAGs have been  
407 criticized for their incompleteness and potentially high prevalence of misassemblies<sup>30</sup>,  
408 we note that i) the overall completeness of our MAGs was rather high (90% median  
409 completeness), ii) complete genomes are not required for accurate taxonomic  
410 profiling<sup>31</sup>, iii) the prevalence of misassemblies among MAGs is likely quite low when  
411 using state-of-the-art assembly and binning approaches<sup>32</sup>. Still, researchers who may  
412 utilize this set of MAGs should use caution when analyzing individual SNPs, plasmids,  
413 genomic islands, or other potentially missing or misassembled genomic features<sup>33</sup>.

414         We investigated the distribution of our MAGs across environment and host  
415 biomes to elucidate the diversity of host-microbe symbiosis in the vertebrate gut.

416 Microbe-host symbiosis spans the continuum from free-living microbes that can simply  
417 survive passage through the host gut, to obligate symbioses<sup>34</sup>. Therefore, MAGs  
418 enriched in the environment versus the host would indicate a weak association, while  
419 the opposite enrichment would suggest a more obligate symbiosis. We provide  
420 evidence of host specificity for the majority of SGBs, while a few Proteobacteria and  
421 Actinobacteria SGBs were environment-enriched. When just considering host-  
422 associated metagenomes, these env-SGBs were generally enriched in non-mammals  
423 (Figures 2, 3, & S8). This is consistent with the hypothesis that mixed-mode  
424 transmission, especially between environmental sources and hosts, is more  
425 commonplace in non-mammalian gut microbiome community assembly versus in  
426 mammals<sup>35</sup>.

427 Our trait-based analysis of SGBs supports the notion that host-enriched taxa are  
428 adapted for a symbiotic lifestyle, while environment-enriched taxa are adapted for a  
429 free-living or facultative symbiosis lifestyle (Figure 3). For instance, anaerobes  
430 comprised almost all host-enriched SGBs, while environment-enriched SGBs were  
431 aerobes or facultative anaerobes and generally motile, which could be highly beneficial  
432 for transmission between the environment and gut biomes. Indeed, a recent directed  
433 evolution experiment showed that selecting for inter-host migration can generate  
434 bacterial strains with increased motility<sup>36</sup>, and a trait-based study of the human infant  
435 gut microbiome showed that later stages of succession are dominated by taxa adapted  
436 to the anoxic gut<sup>37</sup>.

437 By assessing SGB enrichment in Mammalia versus non-Mammalia metagenomes,  
438 we elucidated the specificity of host-microbe symbioses in the gut across large  
439 evolutionary distances. More SGBs were enriched in mammals versus non-mammals  
440 (Figures 2 & S8), as we observed in our previous 16S rRNA assessment of these  
441 vertebrate clades<sup>12</sup>. Few traits differed among SGBs enriched in either biome (Figure  
442 S8), which may indicate that the traits assessed are similarly required for adaptation to  
443 each host clade, even at this coarse evolutionary scale.

444 Vertebrates both play a critical role in the spread of antimicrobial resistance and  
445 also have been sources of novel antibiotics and other natural products<sup>27,38</sup>. We  
446 investigated BGC and AMR diversity in our MAG collection and observed a high  
447 diversity of BGC products, but very few of the BGCs clustered into families with  
448 experimentally characterized BGCs from the MIBiG database (Figures S9 & S10). This  
449 contrasts with findings that only ~10% of BGCs in the human microbiome are  
450 uncharacterized<sup>39</sup>, which is likely due to the limited study of natural products in the gut  
451 microbiome of non-human vertebrates<sup>40,41</sup>. We found NRPS-producing BGCs to be  
452 prevalent among the Firmicutes SGBs, which is similar to a recent assessment of 501  
453 genomes from rumen isolates in which thousands of BGCs were identified<sup>42</sup>. Still, RiPPs  
454 were most prevalent across all vertebrate clades, which expands upon observations of  
455 high prevalence of this BGC class in the gut microbiome of humans<sup>39</sup> (Figures 2 & 4).



456 By combining our AMR marker screen with our SGB biome enrichment analysis,  
457 we were able to characterize how AMR is associated with varying degrees of symbiosis  
458 (Figure S9), which is important for understanding AMR reservoirs<sup>27,43</sup>. Our findings  
459 indicate that the AMR reservoir may be greater for free-living and facultatively symbiotic  
460 taxa relative to microbes with stronger host associations (Figure S9). Indeed, some of  
461 the most abundant AMR markers were associated with metal resistance (*e.g.*, *ruvB*,  
462 *tupC*, and *arsT*), which may reflect a lifestyle in which the microbe is exposed to  
463 environmental sources of metals<sup>44,45</sup>.

464 While MAGs provide a powerful means of investigating species and strain-level  
465 diversity within the vertebrate gut microbiome, the approach is limited to only relatively  
466 abundant taxa with enough coverage to reach adequate assembly contiguity<sup>46</sup>. Our  
467 gene-based assembly approach allowed us to greatly expand the known gene  
468 catalogue of the vertebrate gut microbiome beyond just the abundant taxa, with a total  
469 of >150 million non-redundant coding sequences generated, comprising 88 bacterial  
470 and 11 archaeal phyla (Figure 5). In comparison, recent large-scale metagenome  
471 assemblies of the gut microbiome from chickens, pigs, rats, and dogs have generated  
472 7.7, 9.04, 7.7, 5.1, and 1.25 million non-redundant coding sequences,  
473 respectively<sup>8,25,47,48</sup>. It is also illustrative to consider that a recent large-scale  
474 metagenome assembly of cattle rumen metagenomes generated 69,678 non-redundant  
475 genes involved in carbohydrate metabolism<sup>9</sup>, while our gene collection comprised  
476 substantially more CAZy-annotated gene clusters ( $n = 87,573$ ), even after collapsing at  
477 50% sequence identity. The increased mappability that we achieved across all 5  
478 vertebrate clades when incorporating our gene catalogue in our functional metagenome  
479 profiling pipeline demonstrates how our gene collection will likely aid future vertebrate  
480 gut metagenome studies (Figure S17).

481 Our assessment of gene cluster abundances in metagenomes from environment  
482 and host-associated biomes illuminates how microbiome functioning and taxonomy is  
483 distributed across the free-living to obligate symbiont spectrum. Most notably, nearly all  
484 prominent functional groups were enriched in both the environment and host-associated  
485 biomes, but the specific gene clusters belonged to different taxonomic groups in each  
486 biome (Figure 6). For instance, almost all abundant CAZy families were enriched in both  
487 the environment and host biomes, but the environment was dominated by  
488 Proteobacteria, while Firmicutes, Bacteroidetes, and Actinobacteria gene clusters  
489 comprised most host-enriched CAZy families. This suggests the same coarse-level  
490 functional groups are present across the free-living to obligate microbe-vertebrate  
491 symbiosis lifestyles, but coarse-level taxonomy strongly differs across this spectrum.  
492 This pattern largely remained true when we compared enrichment between the  
493 Mammalia and non-mammals, suggesting that taxonomic differences prevail over  
494 functional differences in regards to host specificity, at least over broad-scale vertebrate  
495 evolutionary distances. While comparing function to taxonomy is challenging due to

496 differing levels of resolution, we do not believe that our findings are simply due to using  
497 functional groupings that are coarser than taxonomy, given that i) we assessed multiple  
498 functional grouping (COG, KEGG, and CAZy), which all showed similar patterns, even  
499 though they differ in functional resolution, and ii) we assessed taxonomy at the very  
500 coarse phylum level but still found stark taxonomic differences across biomes.

501 In conclusion, our large-scale metagenome assembly of both MAGs and coding  
502 sequences from a diverse collection of vertebrates substantially expands the known  
503 taxonomic and functional diversity of the vertebrate gut microbiome. We have  
504 demonstrated that both taxonomic and functional metagenome profiling of the  
505 vertebrate gut is improved by our MAG and gene catalogues, which will aid future  
506 investigations of the vertebrate gut microbiome. Moreover, our collection can help guide  
507 natural product discovery and bioprospecting of novel carbohydrate-active enzymes,  
508 along with modeling AMR transmission among reservoirs. By characterizing the  
509 distribution of MAGs and microbial genes across environment and host biomes, we  
510 gained insight into how taxonomy and function differ along the free-living to obligate  
511 symbiosis lifestyle spectrum. We must note that our metagenome assembly dataset is  
512 biased toward certain animal clades, which likely impacts these findings. As  
513 metagenome assembly becomes more commonplace for studying the vertebrate gut  
514 microbiome, bias toward certain vertebrates (*e.g.*, humans) will decrease, and thus  
515 allow for a more comprehensive reassessment of our findings.

## 516 **Acknowledgements**

517 We thank Nadine Ziemert for helpful discussions in regards to bioinformatic  
518 approaches for secondary metabolite detection and analysis. This work was supported  
519 by the Department of Microbiome Science at the Max Planck Institute for  
520 Developmental Biology. This study was supported by the Austrian Science Fund (FWF)  
521 research projects P23900 granted to Andreas H. Farnleitner and P22032 granted to  
522 Georg H. Reischer. Further support came from the Science Call 2015 "Resource und  
523 Lebensgrundlage Wasser" Project SC15-016 funded by the Niederösterreichische  
524 Forschungs- und Bildungsgesellschaft (NFB).

525 We would like to thank the following collaborators for their huge efforts in sample  
526 and data collection: Mario Baldi, School of Veterinary Medicine, Universidad Nacional  
527 de Costa Rica; Wolfgang Vogl and Frank Radon, Konrad Lorenz Institute of Ethology  
528 and Biological Station Illmitz; Endre Sós and Viktor Molnár, Budapest Zoo; Ulrike  
529 Streicher, Conservation and Wildlife Management Consultant, Vietnam; Katharina Mahr,  
530 Konrad Lorenz Institute of Ethology, University of Veterinary Medicine Vienna and  
531 Flinders University Adelaide, South Australia; Peggy Rismiller, Pelican Lagoon  
532 Research Centre, Australia; Rob Deaville, Institute of Zoology, Zoological Society of  
533 London; Alex Lécu, Muséum National d'Histoire Naturelle and Paris Zoo; Danny  
534 Govender and Emily Lane, South African National Parks, Sanparks; Fritz Reimoser,

535 Research Institute of Wildlife Ecology, University of Veterinary Medicine Vienna; Anna  
536 Küber-Heiss and Team, Pathology, Research Institute of Wildlife Ecology, University of  
537 Veterinary Medicine Vienna; Nikolaus Eisank, Nationalpark Hohe Tauern, Kärnten;  
538 Attila Hettyey and Yoshan Moodley, Konrad Lorenz Institute of Ethology, University of  
539 Veterinary Medicine Vienna; Mansour El-Matbouli and Oskar Schachner, Clinical Unit of  
540 Fish Medicine, University of Veterinary Medicine; Barbara Richter, Institute of Pathology  
541 and Forensic Veterinary Medicine, University of Veterinary Medicine Vienna; Hanna  
542 Vielgrader and Zoovet Team, Schönbrunn Zoo; Reinhard Pichler, Herberstein Zoo. We  
543 explicitly thank the Freek Venter of South African National Parks and the National  
544 Zoological Gardens of South Africa for granting access to their Parks for sample  
545 collection.

#### 546 **Author Contributions**

547 G.H.R., R.E.L., and A.H.F. created the study concept. G.H.R., N.S., C.W., and  
548 G.S. performed the sample collection and metadata compilation. G.H.R., N.S., and S.D.  
549 performed the laboratory work. N.D.Y. and J.C. performed the data analysis. N.D.Y.,  
550 J.C., and R.E.L. wrote the manuscript.

#### 551 **Competing Interest Statement**

552 No conflicts of interest declared.

#### 553 **References**

- 554 1. Nayfach, S., Rodriguez-Mueller, B., Garud, N. & Pollard, K. S. An integrated  
555 metagenomics pipeline for strain profiling reveals novel patterns of bacterial  
556 transmission and biogeography. *Genome Res.* **26**, 1612–1625 (2016).
- 557 2. Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome  
558 research. *BMC Biol.* **17**, 48 (2019).
- 559 3. Wang, W.-L. *et al.* Application of metagenomics in the human gut microbiome.  
560 *World J. Gastroenterol.* **21**, 803–814 (2015).
- 561 4. Zou, Y. *et al.* 1,520 reference genomes from cultivated human gut bacteria enable  
562 functional microbiome analyses. *Nat. Biotechnol.* **37**, 179–185 (2019).
- 563 5. Forster, S. C. *et al.* A human gut bacterial genome and culture collection for  
564 improved metagenomic analyses. *Nat. Biotechnol.* **37**, 186–192 (2019).
- 565 6. Mukherjee, S. *et al.* 1,003 reference genomes of bacterial and archaeal isolates  
566 expand coverage of the tree of life. *Nat. Biotechnol.* **35**, 676–683 (2017).
- 567 7. Almeida, A. *et al.* A unified sequence catalogue of over 280,000 genomes obtained  
568 from the human gut microbiome. *bioRxiv* 762682 (2019) doi:10.1101/762682.
- 569 8. Huang, P. *et al.* The chicken gut metagenome and the modulatory effects of plant-  
570 derived benzylisoquinoline alkaloids. *Microbiome* **6**, 211 (2018).
- 571 9. Stewart, R. D. *et al.* Assembly of 913 microbial genomes from metagenomic  
572 sequencing of the cow rumen. *Nat. Commun.* **9**, 870 (2018).

- 573 10. Riiser, E. S. *et al.* Switching on the light: using metagenomic shotgun sequencing  
574 to characterize the intestinal microbiome of Atlantic cod. *Environ. Microbiol.* **21**,  
575 2576–2594 (2019).
- 576 11. Gibson, K. M. *et al.* Gut microbiome differences between wild and captive black  
577 rhinoceros - implications for rhino health. *Sci. Rep.* **9**, 7570 (2019).
- 578 12. Youngblut, N. D. *et al.* Host diet and evolutionary history explain different aspects  
579 of gut microbiome diversity among vertebrate clades. *Nat. Commun.* **10**, 2200  
580 (2019).
- 581 13. Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time.  
582 *Nat. Commun.* **9**, 2542 (2018).
- 583 14. Steinegger, M., Mirdita, M. & Söding, J. Protein-level assembly increases protein  
584 sequence recovery from metagenomic samples manyfold. *Nat. Methods* **16**, 603–  
585 606 (2019).
- 586 15. Karasov, T. L. *et al.* Arabidopsis thaliana and Pseudomonas Pathogens Exhibit  
587 Stable Associations over Evolutionary Timescales. *Cell Host Microbe* **24**, 168–  
588 179.e4 (2018).
- 589 16. Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: estimating species  
590 abundance in metagenomics data. *PeerJ Comput. Sci.* **3**, e104 (2017).
- 591 17. de la Cuesta-Zuluaga, J., Ley, R. E. & Youngblut, N. D. Struo: a pipeline for  
592 building custom databases for common metagenome profilers. *Bioinformatics*  
593 (2019) doi:10.1093/bioinformatics/btz899.
- 594 18. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and  
595 metatranscriptomes. *Nat. Methods* **15**, 962–968 (2018).
- 596 19. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new  
597 method for improved phylogenetic and taxonomic placement of microbes. *Nat.*  
598 *Commun.* **4**, 2304 (2013).
- 599 20. Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining  
600 pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).
- 601 21. Hannigan, G. D. *et al.* A deep learning genome-mining strategy for biosynthetic  
602 gene cluster prediction. *Nucleic Acids Res.* **47**, e110 (2019).
- 603 22. Navarro-Muñoz, J. C. *et al.* A computational framework to explore large-scale  
604 biosynthetic diversity. *Nat. Chem. Biol.* **16**, 60–68 (2020).
- 605 23. Breitwieser, F. P., Baker, D. N. & Salzberg, S. L. KrakenUniq: confident and fast  
606 metagenomics classification using unique k-mer counts. *Genome Biol.* **19**, 198  
607 (2018).
- 608 24. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and  
609 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 610 25. Coelho, L. P. *et al.* Similarity of the dog and human gut microbiomes in gene  
611 content and response to diet. *Microbiome* **6**, 72 (2018).
- 612 26. Weimann, A. *et al.* From Genomes to Phenotypes: Traitair, the Microbial Trait  
613 Analyzer. *mSystems* **1**, (2016).
- 614 27. Allen, H. K. *et al.* Call of the wild: antibiotic resistance genes in natural  
615 environments. *Nat. Rev. Microbiol.* **8**, 251–259 (2010).
- 616 28. Rodriguez-R, L. M., Gunturu, S., Tiedje, J. M., Cole, J. R. & Konstantinidis, K. T.  
617 Nonpareil 3: Fast Estimation of Metagenomic Coverage and Sequence Diversity.  
618 *mSystems* **3**, (2018).

- 619 29. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny  
620 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
- 621 30. Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A. M. & Banfield, J. F. Accurate  
622 and complete genomes from metagenomes. *Genome Res.* **30**, 315–333 (2020).
- 623 31. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to  
624 classify genomes with the Genome Taxonomy Database. *Bioinformatics* (2019)  
625 doi:10.1093/bioinformatics/btz848.
- 626 32. Mineeva, O., Rojas-Carulla, M., Ley, R. E., Schölkopf, B. & Youngblut, N. D.  
627 DeepMAseD: Evaluating the quality of metagenomic assemblies. *Bioinformatics*  
628 (2020) doi:10.1093/bioinformatics/btaa124.
- 629 33. Maguire, F. *et al.* Metagenome-Assembled Genome Binning Methods with Short  
630 Reads Disproportionately Fail for Plasmids and Genomic Islands.  
631 2020.03.31.997171 (2020) doi:10.1101/2020.03.31.997171.
- 632 34. Sachs, J. L., Skophammer, R. G. & Regus, J. U. Evolutionary transitions in  
633 bacterial symbiosis. *Proc. Natl. Acad. Sci. U. S. A.* **108 Suppl 2**, 10800–10807  
634 (2011).
- 635 35. Ebert, D. The Epidemiology and Evolution of Symbionts with Mixed-Mode  
636 Transmission. *Annu. Rev. Ecol. Evol. Syst.* **44**, 623–643 (2013).
- 637 36. Robinson, C. D. *et al.* Experimental bacterial adaptation to the zebrafish gut reveals  
638 a primary role for immigration. *PLoS Biol.* **16**, e2006893 (2018).
- 639 37. Guittar, J., Shade, A. & Litchman, E. Trait-based community assembly and  
640 succession of the infant gut microbiome. *Nat. Commun.* **10**, 512 (2019).
- 641 38. Adnani, N., Rajsiki, S. R. & Bugni, T. S. Symbiosis-inspired approaches to antibiotic  
642 discovery. *Nat. Prod. Rep.* **34**, 784–814 (2017).
- 643 39. Donia, M. S. *et al.* A systematic analysis of biosynthetic gene clusters in the human  
644 microbiome reveals a common family of antibiotics. *Cell* **158**, 1402–1414 (2014).
- 645 40. Donia, M. S. & Fischbach, M. A. Small molecules from the human microbiota.  
646 *Science* **349**, 1254766 (2015).
- 647 41. Milshteyn, A., Colosimo, D. A. & Brady, S. F. Accessing Bioactive Natural Products  
648 from the Human Microbiome. *Cell Host Microbe* **23**, 725–736 (2018).
- 649 42. Seshadri, R. *et al.* Cultivation and sequencing of rumen microbiome members from  
650 the Hungate1000 Collection. *Nat. Biotechnol.* **36**, 359–367 (2018).
- 651 43. von Wintersdorff, C. J. H. *et al.* Dissemination of Antimicrobial Resistance in  
652 Microbial Ecosystems through Horizontal Gene Transfer. *Front. Microbiol.* **7**, 173  
653 (2016).
- 654 44. Pal, C., Bengtsson-Palme, J., Kristiansson, E. & Larsson, D. G. J. Co-occurrence of  
655 resistance genes to antibiotics, biocides and metals reveals novel insights into their  
656 co-selection potential. *BMC Genomics* **16**, 964 (2015).
- 657 45. Ben Fekih, I. *et al.* Distribution of Arsenic Resistance Genes in Prokaryotes. *Front.*  
658 *Microbiol.* **9**, 2473 (2018).
- 659 46. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes  
660 substantially expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).
- 661 47. Xiao, L. *et al.* A reference gene catalogue of the pig gut microbiome. *Nat Microbiol*  
662 **1**, 16161 (2016).
- 663 48. Pan, H. *et al.* A gene catalogue of the Sprague-Dawley rat gut metagenome.  
664 *Gigascience* **7**, (2018).