

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

3 Nicholas D. Youngblut^{*,1}, Jacobo de la Cuesta-Zuluaga¹, Georg H. Reischer^{2,3}, Silke
4 Dauser¹, Nathalie Schuster², Chris Walzer⁴, Gabrielle Stalder⁴, Andreas H.
5 Farnleitner^{2,3,5}, Ruth E. Ley¹

6 ¹Department of Microbiome Science, Max Planck Institute for Developmental Biology, Max Planck Ring 5,
7 72076 Tübingen, Germany

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

11 ³ICC Interuniversity Cooperation Centre Water & Health, 1160 Vienna, Austria

12 ⁴Research Institute of Wildlife Ecology, University of Veterinary Medicine, Vienna, Austria.

13 ⁵Research Division Water Quality and Health, Karl Landsteiner University for Health Sciences, 3500
14 Krems an der Donau, Austria

15 * Corresponding author: Nicholas Youngblut

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

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

18 **Abstract**

19 Large-scale metagenome assemblies of human microbiomes have produced a
20 vast catalogue of previously unseen microbial genomes; however, comparatively few
21 microbial genomes derive from other vertebrates. Here, we generated 4374
22 metagenome assembled genomes (MAGs) from gut samples of 180 predominantly wild
23 animal species representing 5 classes. Combined with existing datasets, we produced
24 5596 non-redundant, quality MAGs and 1522 species-level genome bins (SGBs). Most
25 SGBs were novel at the species, genus, or family levels, and the majority were enriched
26 in host versus environment metagenomes. Many traits distinguished SGBs enriched in
27 host or environmental biomes, including the number of antimicrobial resistance genes.
28 We identified 1986 diverse and largely novel biosynthetic gene clusters. Gene-based
29 assembly revealed tremendous gene diversity, much of it host or environment specific.
30 Our MAG and gene datasets greatly expand the microbial genome repertoire and
31 provide a broad view of microbial adaptations to life within a living host.

32 **Introduction**

33 The vertebrate gut microbiome comprises a vast amount of genetic diversity, yet
34 even for the most well-studied species such as humans, the number of microbial
35 species lacking a reference genome was recently estimated to be 40-50%¹. Uncovering
36 this “microbial dark matter” is essential to understanding the roles of individual
37 microbes, their intra- and inter-species diversity within and across host populations, and
38 how each microbe interacts with each other and the host to mediate host physiology in
39 a myriad number of ways². On a more applied level, characterizing novel gut microbial
40 diversity aids in bioprospecting of novel bioactive natural products, catalytic and
41 carbohydrate-binding enzymes, probiotics, etc., along with aiding in the discovery and
42 tracking of novel pathogens and antimicrobial resistance (AMR)³.

43 Recent advances in culturomic approaches have generated thousands of novel
44 microbial genomes⁴⁻⁶, but the throughput is currently far outpaced by metagenome
45 assembly approaches⁷. However, such large-scale metagenome assembly-based
46 approaches have not been as extensively applied to most non-human vertebrates. The
47 low amount of metagenome reads classified in some recent studies of the rhinoceros,
48 chicken, cod, and cow gut/rumen microbiome suggests that databases lack much of the
49 genomic diversity in less-studied vertebrates⁸⁻¹¹. Indeed, the limited number of studies
50 incorporating metagenome assembly hint at the extensive amounts of as-of-yet novel
51 microbial diversity across the >66,000 vertebrate species on our planet.

52 Here, we developed an extensive metagenome assembly pipeline and applied it
53 to a multi-species dataset of microbiome diversity across vertebrate species comprising

54 5 classes: Mammalia, Aves, Reptilia, Amphibia, and Actinopterygii, with >80% of
55 samples obtained from wild individuals¹² combined with data from
56 14 published animal gut metagenomes. Moreover, we also applied a recently developed
57 gene-based metagenome assembly pipeline to the entire dataset in order to obtain
58 gene-level diversity for rarer taxa that would otherwise be missed by genome-base
59 assembly^{13,14}. Our assembly approaches yielded a great deal of novel genetic diversity,
60 which we found to be largely enriched in animals versus the environment, and to some
61 degree, enriched in particular animal clades.

62 **Methods**

63 *Sample collection*

64 Sample collection was as described in Youngblut and colleagues¹². Table S1 shows
65 the dates, locations, and additional metadata of all samples collected. All fecal samples
66 were collected in sterile sampling vials, transported to a laboratory and frozen within 8
67 hours. DNA extraction was performed with the PowerSoil DNA Isolation Kit (MoBio
68 Laboratories, Carlsbad, USA).

69 *“multi-species” vertebrate gut metagenomes*

70 Metagenome libraries were prepared as described by Karasov and colleagues¹⁵.
71 Briefly, 1 ng of input gDNA was used for Nextera Tn5 tagmentation. A BluePippin was
72 used to restrict fragment sizes to 400-700 bp. Barcoded samples were pooled and
73 sequenced on an Illumina HiSeq3000 with 2x150 paired-end sequencing. Read quality
74 control (QC) is described in the Supplemental Methods.

75 Post-QC reads were taxonomically profiled with Kraken2 and Bracken v.2.2¹⁶
76 against the Struo-generated GTDB-r89 Kraken2 and Bracken databases¹⁷. HUMAnN2
77 v.0.11.2¹⁸ was used to profile genes and pathways against the Struo-generated
78 HUMAnN2 database created from GTDB-r89.

79 *Publicly available animal gut metagenomes*

80 Published animal gut metagenome reads were downloaded from the Sequence
81 Read Archive (SRA) between May and August of 2019. Table S2 lists all included
82 studies. We selected studies with Illumina paired-end metagenomes from gut contents
83 or feces. MGnify samples were downloaded from the SRA in Oct 2019 (Table S3). Read
84 quality control is described in the Supplemental Methods.

85 *Metagenome assembly of genomes pipeline*

86 Assemblies were performed on a per-sample basis, with reads subsampled via
87 seqtk v.1.3 to ≤ 20 million read pairs. The details of the assembly pipeline are described
88 in the Supplemental Methods.

89 A multi-locus phylogeny of all SGB representatives was inferred with PhyloPhlan
90 v.0.41¹⁹. Secondary metabolites were identified with AntiSMASH v.5.1.1²⁰ and
91 DeepBGC v.0.1.18²¹ and then characterized with BiGSCAPE²². Abriicate was used to
92 identify antimicrobial resistance genes. We used Krakenuniq v.0.5.8²³ for estimating
93 abundance of MAGs in metagenome samples. Details can be found in the
94 Supplemental Methods.

95 *Metagenome assembly of genes pipeline*

96 Assemblies performed on a per-sample basis, with reads subsampled via seqtk
97 v.1.3 to ≤ 20 million pairs. We used PLASS v.2.c7e35¹⁴ and Linclust (mmseqs
98 v.10.6d92c)¹³ to assemble and cluster contigs. A full description is in the Supplemental
99 Methods. DESeq2²⁴ was used to estimate enrichment of MAGs and gene clusters in
100 metagenomes from host and environment biomes.

101 *Data availability*

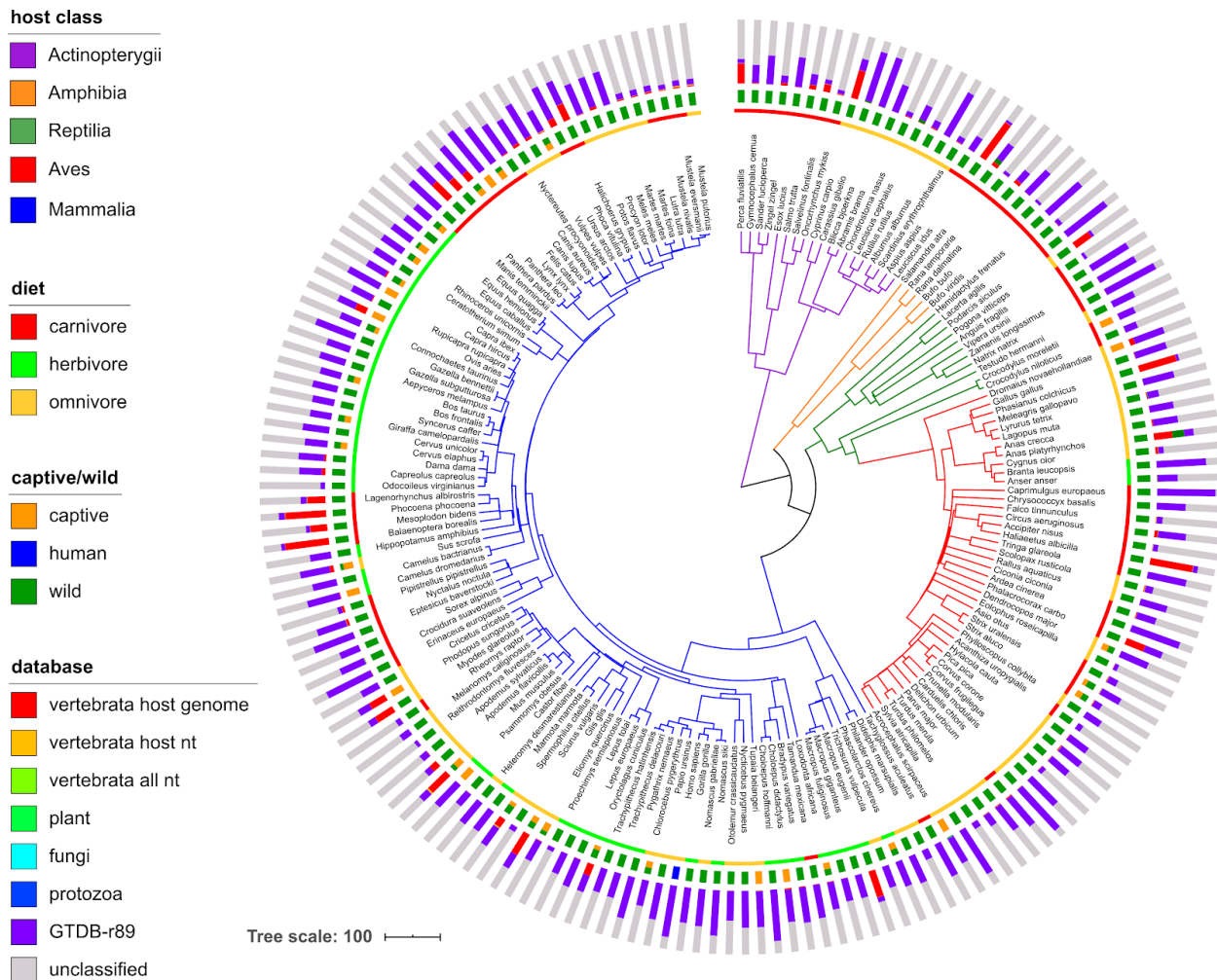
102 The raw sequence data are available from the European Nucleotide Archive
103 (ENA) under the study accession number PRJEB38078. Fasta files for the 5596
104 non-redundant MAGs, 1522 SGBs, and gene clusters (50, 90, and 100% sequence
105 identity clustering) can be found at
106 http://ftp.tue.mpg.de/ebio/projects/animal_gut_metagenome_assembly/. Code used for
107 processing the data can be found at
108 https://github.com/leylabmpi/animal_gut_metagenome_assembly.

109 **Results**

110 *Animal gut metagenomes from a highly diverse collection of animals*

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

117 Our read-quality control pipeline included stringent filtering of host reads; some
 118 samples contained high amounts of reads mapping to vertebrate genomes (up to 74%;
 119 $6 \pm 17\%$ s.d.; Figure 1). Gut content samples contained a significantly higher amount of
 120 host reads ($13.5 \pm 21.6\%$ s.d.) versus feces metagenomes ($4.7 \pm 12.7\%$ s.d.; Wilcox, P
 121 $< 1.8e-7$; Table S1). We mapped all remaining reads to a custom comprehensive
 122 Kraken2 database built from the GTDB (Release 89). Still, many samples had a low
 123 percentage of mapped reads (43 ± 22 s.d.; Figure 1), with 29% of the samples having
 124 $<20\%$ mapped reads.



125 **Figure 1.** A large percentage of unmapped reads, even when using multiple comprehensive metagenome
 126 profiling databases. The dated host species phylogeny was obtained from <http://timetree.org>, with
 127 branches colored by host class. From inner to outer, the data mapped onto the tree is host diet, host
 128 captive/wild status, and the mean number of metagenome reads mapped to various host-specific,
 129 non-microbial, and microbial databases. Note that captive/wild status sometimes differs among individuals
 130 of the same species. The databases are i) a representative of each publicly available genome from the
 131 host species or sister species ("vertebrata host genome"), ii) all entries in the NCBI nt database with

132 taxonomy IDs matching host species ("vertebrata host nt"), iii) as the previous, but all vertebrata
133 sequences included, iv) the kraken2 "plant" database, v) the kraken2 "fungi" database, vi) the kraken2
134 "protozoa" database, vii) a custom bacteria and archaea database created from the Genome Taxonomy
135 Database, Release 89 ("GTDB-r89"). Reads were mapped iteratively to each database in the order
136 shown in the legend (top to bottom), with only unmapped reads included in the next iteration.
137 "unclassified" reads did not map to any database, which were used along with reads mapping to
138 GTDB-r89 for downstream analyses ("microbial + unclassified").

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

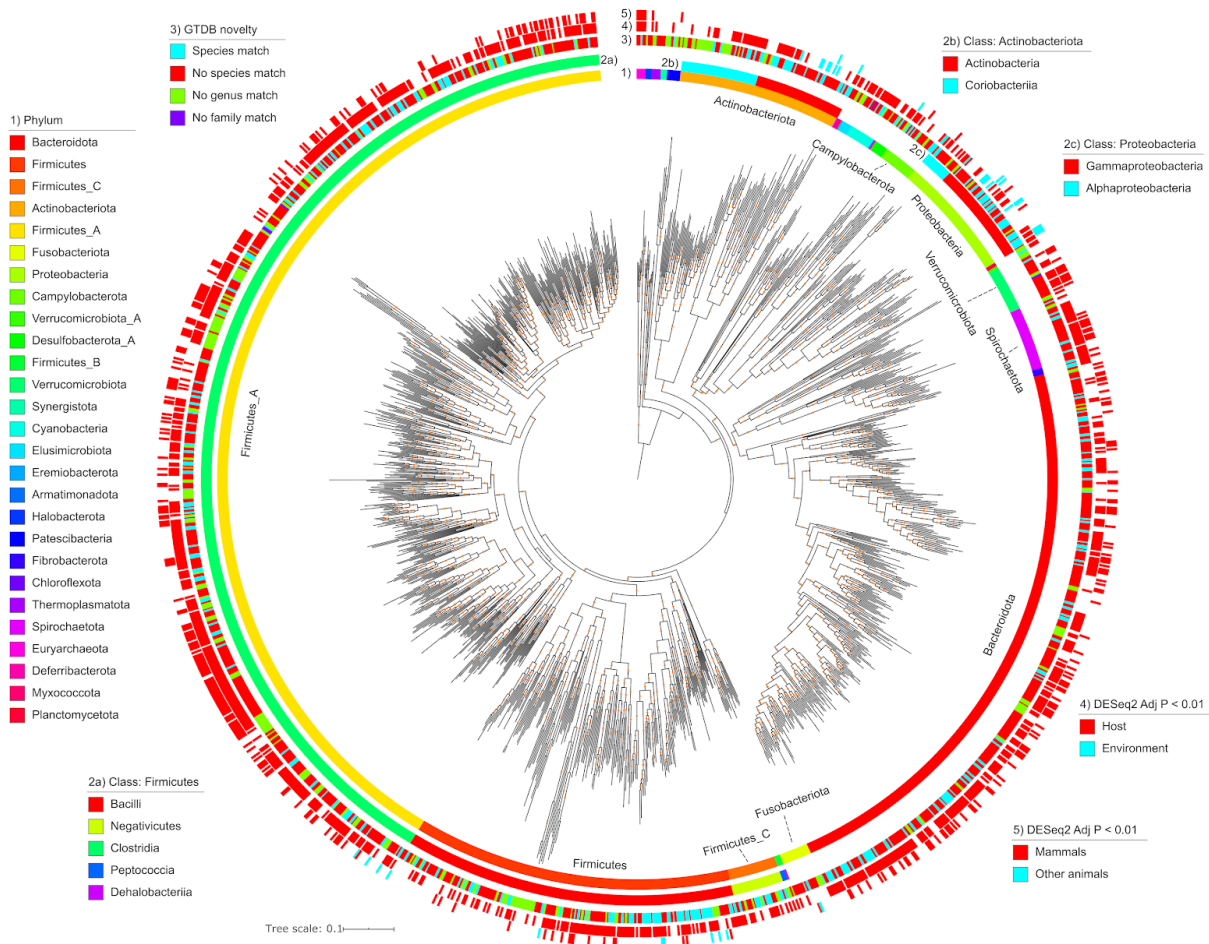
140 Our comprehensive metagenome assembly pipeline generated 4374
141 non-redundant MAGs. After filtering to just "quality" MAGs (see Methods), 296 MAGs
142 remained, with a mean percent completeness and contamination of 84 ± 14 and $1.5 \pm$
143 1.2 s.d., respectively. The MAGs consisted of 11 bacterial and 1 archaeal phylum, as
144 determined via GTDB-Tk²⁵. The majority of MAGs belonged to the classes Clostridia (n
145 = 95; Firmicutes_A phylum) and Bacteroidia ($n = 74$; Bacteroidota phylum; Figure S2).
146 De-replicating MAGs at 95% ANI produced 248 species-level genome bins (SGBs). Of
147 the SGBs, 196 (79%) had <95% ANI to every genome in the GTDB-r89 database, and
148 51 (21%) lacked a genus-level match. These findings indicated that the MAG dataset
149 contained a substantial amount of novel diversity.

150 We expanded our MAG dataset by applying our assembly pipeline to 14
151 publically available animal gut metagenome datasets in which no MAGs have been
152 generated by *de novo* metagenome assembly (Table S2). Our metagenome selection
153 included 554 samples from members of Mammalia (dogs, cats, woodrats, pigs, whales,
154 rhinoceroses, pangolins, and non-human primates), Aves (geese, kakapos, and
155 chickens), and Actinopterygii (cod). We applied our assembly pipeline to each individual
156 dataset and generated a total of 5301 quality MAGs (Figure S3). As with the
157 multi-species metagenome assemblies, MAG quality was high, with a mean
158 completeness and contamination of 85 ± 13 and 1.1 ± 1.1 s.d., respectively. The
159 taxonomic diversity was also quite high, with 2 archaeal and 25 bacterial phyla
160 represented (Figure S3). De-replicating MAGs at 95% ANI produced 1308 SGBs. Of
161 these, 1001 lacked a $\geq 95\%$ ANI match to the GTDB-r89, 216 lacked a genus-level
162 match, and 6 lacked even a family-level match.

163 We combined all quality MAGs and de-replicated at 99.9 and 95% ANI to
164 produce 5596 non-redundant MAGs and 1522 species-level genome bins (SGBs),
165 respectively (Tables S4 & S5). Of the 5596 MAGs, 2773 (50%) had a completeness of
166 $\geq 90\%$. Of the 1522 SGBs, 1184 (78%) lacked a $\geq 95\%$ ANI match to the GTDB-r89, 266
167 (17%) lacked a genus-level match, and 6 lacked a family-level match (Figures 2 & S4).
168 Mapping taxonomic novelty onto a multi-locus phylogeny of all 1522 SGBs revealed that
169 novel taxa were rather dispersed across the phylogeny (Figure 2).

170 We assessed the novelty of our SGBs relative to UHGG, a comprehensive
 171 human gut genome database, and found that only 31% of our SGBs had $\geq 95\%$ ANI to
 172 any of the 4644 UHGG representatives, and this overlap only increased to 34% at a
 173 90% ANI cutoff.

174 Integrating the 1522 SGBs into our custom GTDB Kraken2 database significantly
 175 increased the percent reads mapped (t-test, $P < 0.005$; Figure S5). The percent
 176 increase varied from <1 to 62.8% (mean of 5.3 ± 6.7 s.d.) among animal species but did
 177 not appear biased to just pigs, dogs, or other vertebrate species in the 14 public
 178 metagenome datasets that we incorporated (Figure S6).



179 **Figure 2.** A phylogeny of all 1522 SGBs. From innermost to outermost, the data mapped onto the
 180 phylogeny is: GTDB phylum-level taxonomic classifications, class-level taxonomies for Actinobacteriota,
 181 class-level taxonomies for Firmicutes, class-level taxonomies for Proteobacteria, taxonomic novelty,
 182 significant enrichment in host gut or environmental metagenomes, and significant enrichment in Mammals
 183 or other animals in our multi-species gut metagenome dataset. The phylogeny was inferred from multiple
 184 conserved loci via PhyloPhlAn. Orange dots on the phylogeny denote bootstrap values in the range of 0.7
 185 to 1. The phylogeny is rooted on the last common ancestor of Archaea and Bacteria. The tree scale unit
 186 is substitutions per site.

187 *Enrichment of SGBs among animal clades*

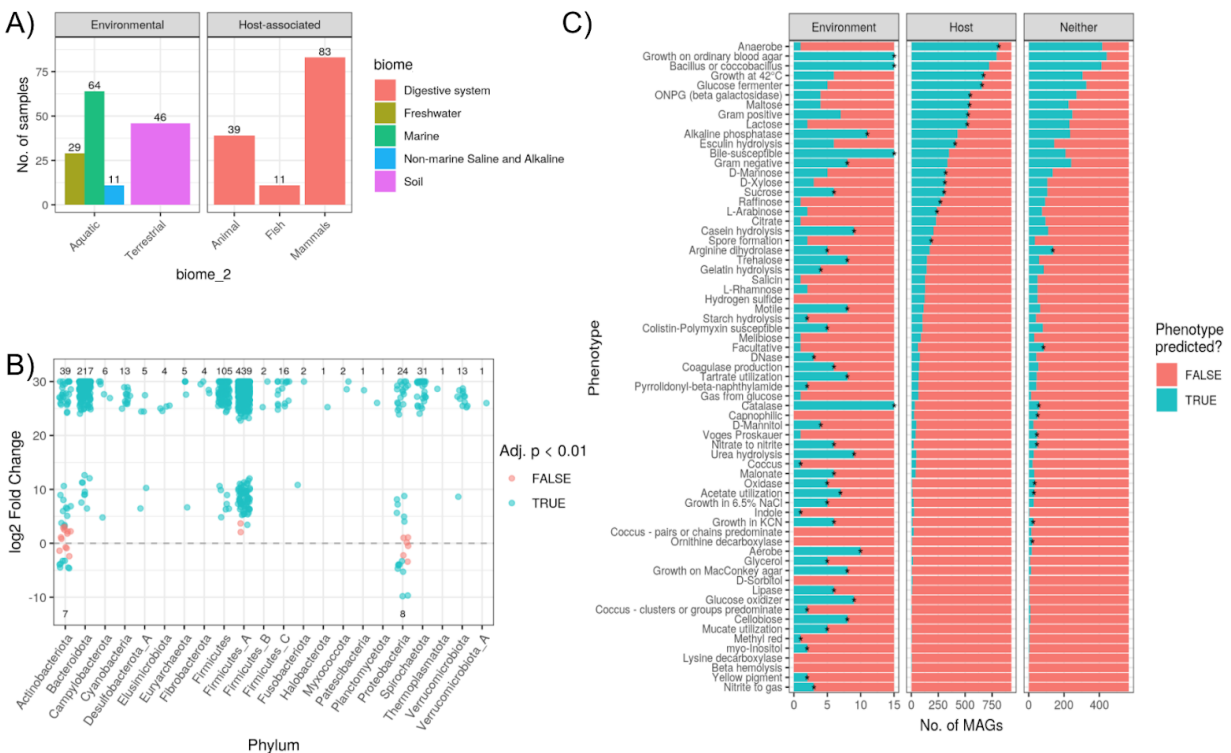
188 While the MAGs generated here derive from animal gut metagenomes, many of
189 these taxa might be transient in the host and actually more prevalent in the
190 environment. We tested this by generating a “host-environment” metagenome dataset
191 comprising 283 samples from 30 BioProjects (17 environmental and 13 host-associated;
192 Figure 3A). We found 932 of the 1522 SGBs (61%) to be significantly enriched in the
193 host metagenomes (DESeq2, *adj. P* < 0.01; Figure 3B). The host-enriched SGBs
194 (host-SGBs) were taxonomically diverse, comprising 22 phyla. In contrast, only 15
195 SGBs (1%) were environment-enriched (env-SGBs), which all belonged to either
196 Actinobacteriota or Proteobacteria (Figure 3B). The only SGBs that were not
197 significantly enriched in either group belonged to Actinobacteriota or Proteobacteria,
198 along with two SGBs from the Firmicutes_A phylum. Mapping these data onto the SGB
199 phylogeny revealed phylogenetic clustering of the environment-enriched SGBs (Figure
200 2).

201 We investigated the traits of the host- and environment-enriched SGBs and
202 found many predicted phenotypes to be more prevalent in one or the other group
203 (Figure 3C). Almost all env-SGBs were aerobes (93%), which may aid in transmission
204 between the environment and host biomes. In contrast, 87% of host-SGBs were
205 anaerobes. Furthermore, all env-SGBs could generate catalase and were bile
206 susceptible, while both phenotypes were sparse in host-SGBs (Figure 3C).
207 Carbohydrate metabolism also differed, with most host-SGBs predicted to consume
208 various tri-, di-, and mono-saccharides. In contrast, env-SGBs were enriched in
209 phenotypes associated with motility, nitrogen metabolism, and breakdown of
210 heterogeneous substrates (*e.g.*, cellobiose metabolism).

211 We also compared SGB enrichment in mammals versus non-mammals in our
212 “multi-species” metagenome dataset and found 361 SGBs (24%) to be significantly
213 enriched in mammals, while 22 (1%) were enriched in non-mammals (DESeq2, *adj. P* <
214 0.01; Figure S7A). Interestingly, 100% of SGBs in the two archaeal phyla (Halobacteria
215 and Euryarchaeota) were enriched in mammals. Also of note, most of the
216 Verrucomicrobiota SGBs (87%) were enriched in mammals. The only 2 phyla with >10%
217 of SGBs enriched in non-mammals were Proteobacteria (29%) and Campylobacteria
218 (25%).

219 In contrast to our assessment of phenotypes distinct to host- or env-SGBs, we
220 did not observe such a distinction of phenotypes among SGBs enriched in Mammalia or
221 non-mammal gut metagenomes (Figure S7B). Certain phenotypes such as anaerobic
222 growth and lactose consumption were more prevalent among mammal species, but they
223 were not found to significantly enriched relative to the null model.

224 Little is known about the distribution of antimicrobial resistance genes in the gut
 225 microbiomes of most vertebrate species²⁶; therefore, we investigated the distribution of
 226 AMR genes among MAGs enriched in the environment versus host biomes. We found a
 227 mean of 35 ± 26 s.d. AMR markers per genome (Figure S8A). The high average was
 228 largely driven by Proteobacteria and Campylobacter genomes, which had a mean of
 229 387 and 161 AMR markers per genome, respectively. The 5 most abundant markers
 230 were *ruvB*, *galE*, *tupC*, *fabL* (*ygaA*), and *arsT* (Figure S8A). The more abundant
 231 markers predominantly belonged to Firmicutes_A, while Proteobacteria comprised
 232 larger fractions of the less abundant markers. Environment-enriched taxa contained
 233 substantially more AMR genes than host-enriched taxa, and the same was true for
 234 non-Mammalia versus Mammalia-enriched taxa (Figure S8B & S8C).

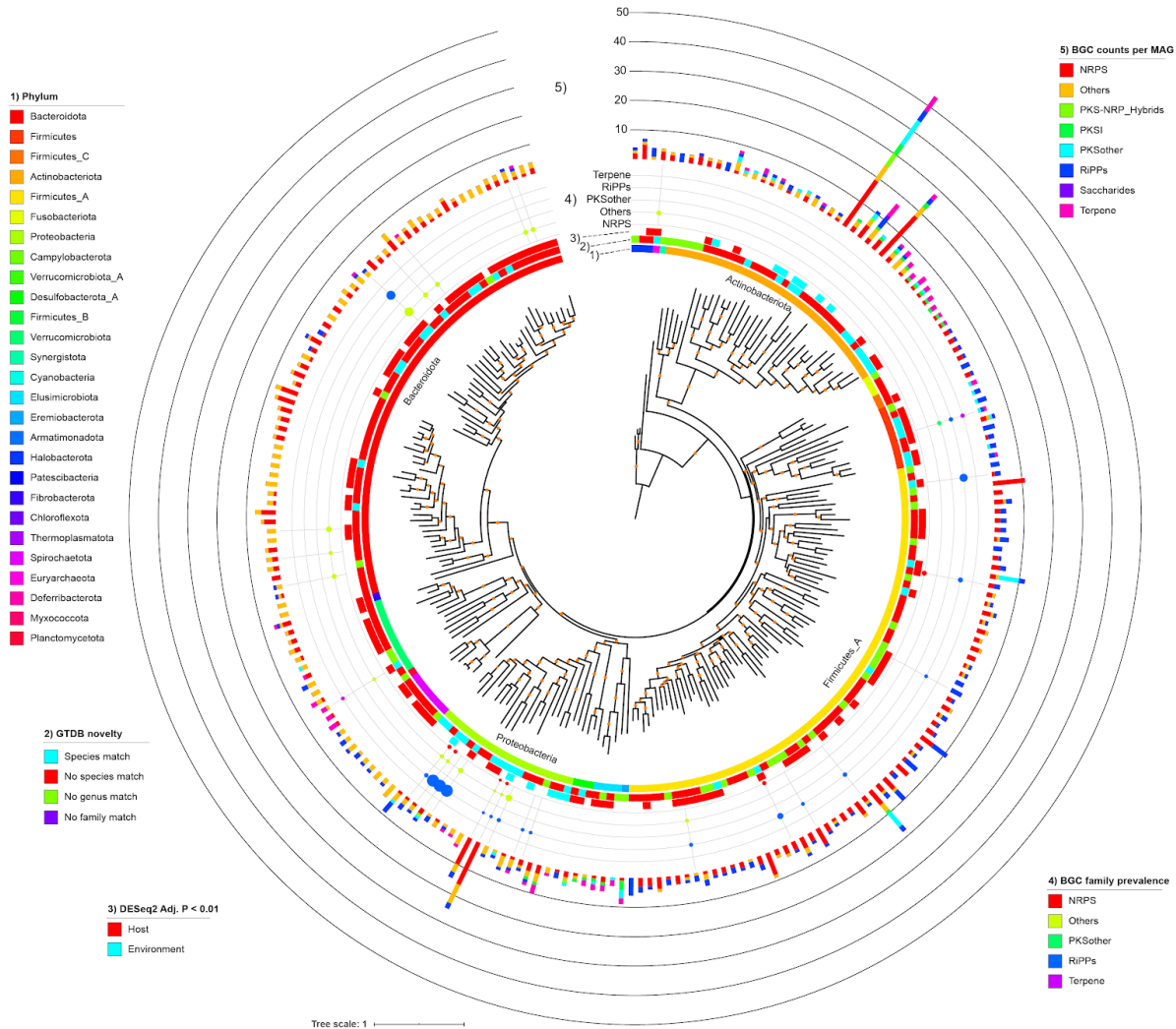


235 **Figure 3.** A) Summary of the number of samples per biome for our multi-environment metagenome
 236 dataset selected from the MGnify database. B) Number of SGBs found to be significantly enriched in host
 237 versus (positive \log_2 fold change; “I2fc”) environmental metagenomes (negative I2fc). Values shown are
 238 the number of MAGs significantly enriched (blue) in either biome or not found to be significant (red). C)
 239 Host- and environment-enriched SGBs have distinct traits. Predicted phenotypes are summarized for the
 240 SGBs significantly enriched in host or environmental metagenomes (DESeq2 *Adj. P* < 0.01) or neither
 241 biome (“Neither” in the x-axis facet). Note the difference in x-axis scale. Asterisks denote phenotypes
 242 significantly more prevalent in SGBs of the particular biome versus a null model of 1000 permutations in
 243 which biome labels were shuffled among SGBs.

244 *MAGs reveal novel secondary metabolite diversity*

245 We identified 1986 biosynthetic gene clusters (BGCs) among all 1522 SGBs. A
246 total of 28 different products were predicted, with the most abundant being
247 non-ribosomal peptide synthetases (NPRS; $n = 473$), sactipeptides ($n = 307$), and
248 arylpolyenes ($n = 291$; Figure S9). BGCs were identified in 2 archaeal and 18 bacterial
249 phyla. MAGs in the Firmicutes_A phylum contained the most BGCs ($n = 764$; 38%),
250 while Bacteroidota and Actinobacteriota phyla possessed 381 (19%) and 272 (14%),
251 respectively (Figure S9). Still, Actinobacteriota SGBs did possess the highest average
252 number of BGCs per genome (16.3), followed by Eremiobacterota (9), Proteobacteria
253 (7.7), and Halobacterota (5.1).

254 Clustering all 1986 BGCs by BiGSCAPE generated 1764 families and 1305
255 clans. BGCs from the MIBiG database only clustered with 8 clans, suggesting a high
256 degree of novelty (Figure S10). Mapping the BGCs on a genome phylogeny of all
257 species containing ≥ 3 BGCs (233 SGBs) revealed that the number of BGCs per
258 genome was somewhat phylogenetically clustered: the five genomes with the most
259 BGCs belonged either to the Actinobacteria or Gammaproteobacteria (Figure 4).
260 Notably, these clades contained a high number of host-SGBs. Of these 233 SGBs, the
261 majority were taxonomically novel, with 62% lacking a species-level match to
262 GTDB-r89, and 18% lacking a genus-level match (Figure 4). To determine which of the
263 BGCs are most prevalent across animal hosts, we quantified the prevalence of each
264 BGC family across our multi-species metagenome dataset and mapped it to the
265 genome phylogeny (Figures 4 & S11). Of the 1543 BGC families, 83 were present in
266 $\geq 25\%$ of the animal metagenomes, with ribosomally synthesized and post-translationally
267 modified peptides (RIPPs) being by far the most prevalent (up to 98% prevalence of
268 individual BGC families) and also found in species from a number of phyla.

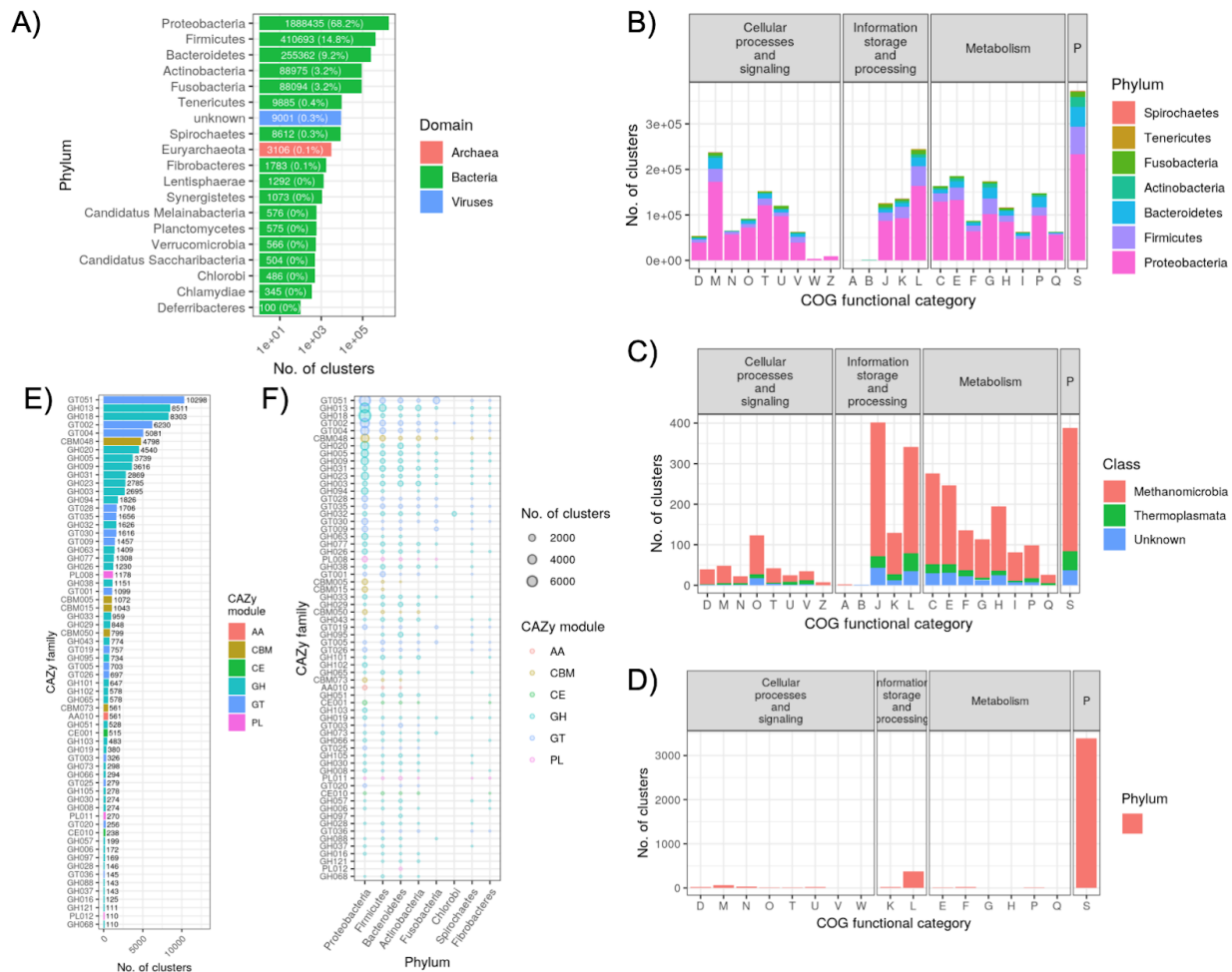


269 **Figure 4.** Phylogeny of all SGBs with ≥ 3 BGCs identified by AntiSMASH. From innermost to outermost,
 270 the data mapped onto the phylogeny is: 1) GTDB phylum-level taxonomic classifications, 2) taxonomic
 271 novelty, 3) significant enrichment in host or environmental metagenomes, 4) the prevalence of BGC
 272 families across the multi-species metagenome dataset, and 5) the number of BGCs identified in the MAG.
 273 Prevalence is the maximum of any BGC family for that BGC type, and only BGC families with a
 274 prevalence of $\geq 25\%$ are shown. The phylogeny is a pruned version of that shown in Figure 2.

275 *Large-scale gene-based metagenome assembly reveals novel diversity*

276 We applied gene-based assembly methods to our combined metagenome
 277 dataset¹⁴, which generated a total of 150,718,125 non-redundant coding sequences
 278 (average length of 179 amino acids). Clustering at 90 and 50% sequence identity
 279 resulted in 140,225,322 and 6,391,861 clusters, respectively. Only 16.9 and 11.3% of
 280 each respective cluster set mapped to the UniRef50 database, indicating that most
 281 coding sequences were novel. The clusters comprised 88 bacterial and 11 archaeal
 282 phyla; 80 of which were represented by <100 clusters, and 60 lacking a cultured

283 representative. Proteobacteria (mostly Gammaproteobacteria), Firmicutes, and
 284 Bacteroidetes made up 92.2% of all clusters (Figure 5A). The proportion of clusters
 285 belonging to each COG functional category was largely the same for the more abundant
 286 bacterial phyla (Figure 5B), while more variation was seen among Euryarchaeota
 287 (Figure 5C). The dominant 7 phyla showed substantial variation in the number of
 288 clusters associated with various KEGG pathway categories (Figure S13). For instance,
 289 a high proportion of Fusobacteria and Tenericutes clusters were associated with the
 290 “nucleotide metabolism”, “replication and repair”, and “translation” categories. A total of
 291 87,573 clusters were annotated as CAZy families, with GT51, GH13, GH18, GT02, and
 292 GT04 representing 48% of all CAZy-annotated clusters (Figure 5E). Of the 12 phyla with
 293 the most CAZy family clusters, there were substantial differences in proportions of
 294 clusters falling into each family (Figure 5F).



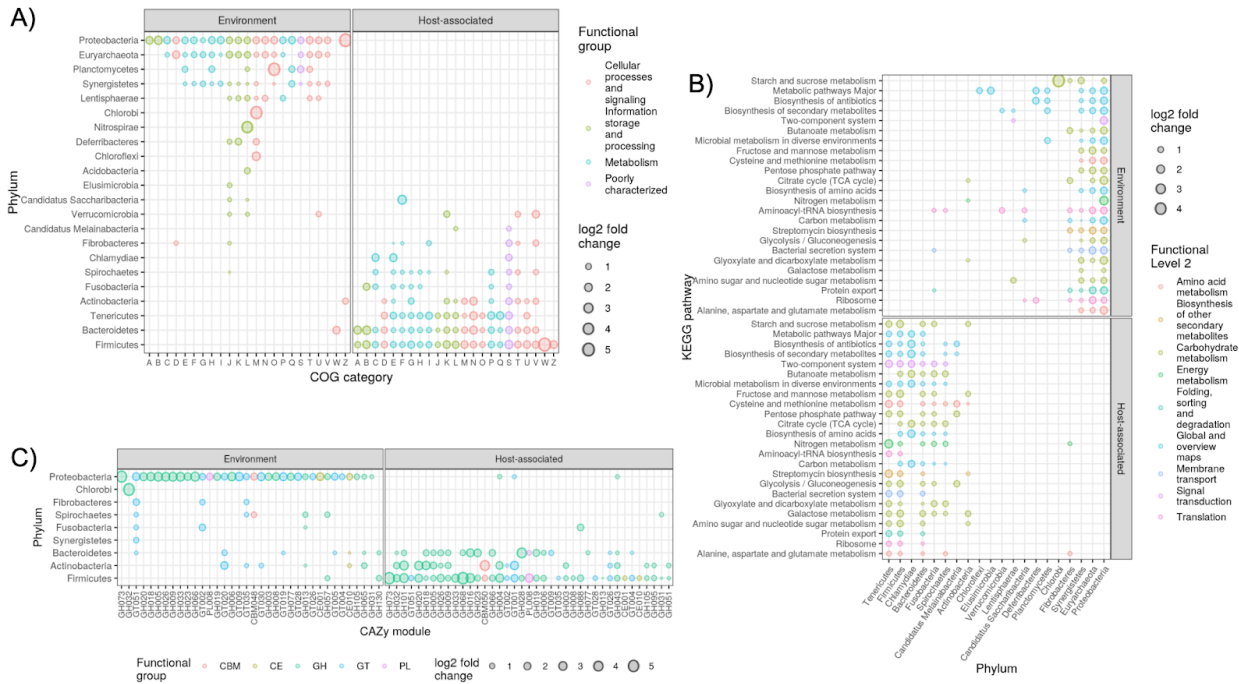
295 **Figure 5.** A summary of the 50% sequence identity clusters generated from the gene-based metagenome
 296 assembly of the combined dataset. A) The total number of gene clusters per phylum. For clarity, only
 297 phyla with ≥ 100 clusters are shown. Labels on each bar list the number of clusters (and percent of the
 298 total). B) The number of bacterial gene clusters per phylum and COG category. The “P” facet label refers

299 to “poorly characterized”. C) The number of archaeal gene clusters per class (all belonging to
300 Euryarchaeota) and COG category. D) The number of viral gene clusters per COG category. E) The
301 number clusters annotated as each CAZy family. For clarity, only phyla with ≥ 100 clusters are shown.
302 Labels next to each bar denote the number of clusters. F) The number of clusters per CAZy family,
303 broken down by phylum. CAZy families and phyla are ordered by most to least number of clusters. For
304 clarity, only CAZy families and phyla with ≥ 100 total clusters are shown.

305 *Biome enrichment of gene clusters from specific phyla*

306 We mapped reads from our host-environment metagenome dataset to each
307 cluster and used DESeq2 to identify those significantly enriched (*adj. P* < 1e-5) in each
308 biome. Most strikingly, the same functional groups were enriched in both biomes,
309 regardless of the grouping (*i.e.*, COG functional category, KEGG pathway, or CAZy
310 family); however, the gene clusters belonged to different microbial phyla (Figure 6;
311 Supplemental Results). For instance, nearly all COG categories for gene clusters
312 belonging to Proteobacteria were environment-enriched, while the same COG
313 categories for clusters belonging to Firmicutes and Bacteroidetes were host-enriched. In
314 contrast, functional groups of certain phyla were enriched in one biome, while different
315 groups were enriched in the other, indicating within-phylum differences in functional
316 content and habitat distributions. For instance, Fusobacteria KEGG pathways were
317 predominantly host-enriched, but protein export, bacteria secretion system, and
318 aminoacyl-tRNA biosynthesis were environment-enriched, indicating that these 3
319 pathways were more predominant in environment-enriched members of Fusobacteria
320 (Figure 6B). Overall, these results suggest that both biomes select for these same
321 microbial functions, but the microbes involved often differ at coarse taxonomic scales.

322 We also assessed gene cluster enrichment in Mammalia versus non-Mammalia
323 and found fewer significantly enriched features, which may be due to the smaller
324 metagenome sample size or less pronounced partitioning of functional groups among
325 biomes (Figure S14; Supplemental Results). Still, we again observed that both biomes
326 enriched for the same microbial functions, but these belonged to different coarse
327 taxonomic groups.



328 **Figure 6.** Enrichment of gene clusters grouped by phylum and A) COG category B) KEGG pathway or C)
 329 CAZy family. Only groupings significantly enriched in abundance (DESeq2, *adj. P* < 1e-5) in either biome
 330 are shown. Only gene clusters observed in at least 25% of the metagenomes were included. For clarity,
 331 only KEGG pathways enriched in >7 phyla are shown, and only CAZy families enriched in >1 phylum are
 332 shown. Note that the axes are flipped in B) relative to A) and C).

333 *Functional metagenome profiling benefits from our gene catalogue*

334 Lastly, We created a custom gene-level metagenome profiling database for the
 335 HUMAnN2 pipeline by merging our coding sequence catalogue with our previously
 336 constructed custom GTDB-r89 database for HUMAnN2¹⁷. We mapped our multi-species
 337 metagenomes to each database via the HUMAnN2 pipeline and compared the percent
 338 reads mapped. Due to the constraint of HUMAnN2 that all references must have a
 339 UniRef ID, we could only use 11.3% (*n* = 722 795) of our gene clusters. Still, we found
 340 that including these clusters increased the mappability by 4 ± 5% s.d. (Figure S15).
 341 Mammalia species benefited the most, but at least one species from each class showed
 342 a mappability increase of >10% (Figure S15B).

343 **Discussion**

344 Our multi-species gut metagenome dataset, derived from >80% wild species
 345 from five vertebrate taxonomic classes, greatly helps to expand the breadth of
 346 cross-species gut metagenome comparisons (Figure 1). By assembling the
 347 metagenomes of our multi-species dataset together with 14 other animal gut
 348 metagenome datasets from understudied host species, we have produced an extensive

349 MAG collection that includes 1184, 266, and 6 genomes from novel species, genera,
350 and families, respectively (Figures 2 & S4). Moreover, we found little overlap (31%)
351 between our MAG collection and the extensive human microbiome genome catalogue
352 comprising the UHGG, which underscores its taxonomic novelty. Our MAG collection,
353 once combined with the GTDB²⁷, improved our ability to classify reads in our
354 multi-species metagenome dataset (Figure S5), which is critical for accurately
355 assessing gut microbiome diversity across vertebrates.

356 We investigated the distribution of our MAGs across environment and host
357 biomes to elucidate the diversity of host-microbe symbiosis in the vertebrate gut.
358 Microbe-host coevolution spans the continuum from free-living microbes that can simply
359 survive passage through the host gut, to obligate symbioses²⁸. Therefore, MAGs
360 enriched in the environment versus the host would indicate a weak association, while
361 the opposite enrichment would suggest a more obligate symbiosis. We provide
362 evidence of host specificity for the majority of SGBs, while a few Proteobacteria and
363 Actinobacteria SGBs were environment-enriched. When just considering
364 host-associated metagenomes, these env-SGBs were generally enriched in
365 non-mammals (Figures 2, 3, & S7). This is consistent with the hypothesis that
366 mixed-mode transmission, especially between environmental sources and hosts, is
367 more commonplace in non-mammalian gut microbiome community assembly versus in
368 mammals²⁹.

369 Our trait-based analysis of SGBs supports the notion that host-enriched taxa are
370 adapted for a symbiotic lifestyle, while environment-enriched taxa are adapted for a
371 free-living or facultative symbiosis lifestyle (Figure 3). For instance, anaerobes
372 comprised almost all host-enriched SGBs, while environment-enriched SGBs were
373 aerobes or facultative anaerobes and generally motile, which could be highly beneficial
374 for transmission between the environment and gut biomes. Indeed, a recent directed
375 evolution experiment showed that selecting for inter-host migration can generate
376 bacterial strains with increased motility³⁰.

377 By assessing SGB enrichment in Mammalia versus non-mammalian
378 metagenomes, we elucidated the specificity of host-microbe symbioses in the gut
379 across large evolutionary distances. More SGBs were enriched in mammals versus
380 non-mammals (Figures 2 & S7), as we observed in our previous 16S rRNA assessment
381 of these vertebrate clades¹². Few traits differed among SGBs enriched in either biome
382 (Figure S7), which may indicate that the traits assessed are similarly required for
383 adaptation to each host clade, even at this coarse evolutionary scale.

384 Vertebrates both play a critical role in the spread of antimicrobial resistance and
385 also have been sources of novel antibiotics and other natural products^{26,31}. We
386 investigated BGC and AMR diversity in our MAG collection and observed a high
387 diversity of BGC products, but very few of the BGCs clustered into families with

388 experimentally characterized BGCs from the MIBiG database (Figures S8 & S9). This
389 contrasts with findings that only ~10% of BGCs in the human microbiome are
390 uncharacterized³², which likely reflects the limited study of natural products in the gut
391 microbiome of non-human vertebrates^{33,34}. Our findings indicate that the AMR reservoir
392 may be greater for free-living and facultatively symbiotic taxa relative to microbes with
393 stronger host associations. Our findings also indicate that AMR may be more prevalent
394 in the guts of non-mammalian hosts (Figure S8).

395 While MAGs provide a powerful means of investigating species and strain-level
396 diversity within the vertebrate gut microbiome, the approach is limited to only relatively
397 abundant taxa with enough coverage to reach adequate assembly contiguity³⁵. Our
398 gene-based assembly approach allowed us to greatly expand the known gene
399 catalogue of the vertebrate gut microbiome beyond just the abundant taxa, with a total
400 of >150 million non-redundant coding sequences generated, comprising 88 bacterial
401 and 11 archaeal phyla (Figure 5). In comparison, recent large-scale metagenome
402 assemblies of the gut microbiome from chickens, pigs, rats, and dogs have generated
403 7.7, 9.04, 7.7, 5.1, and 1.25 million non-redundant coding sequences, respectively^{8,36–38}.
404 It is also illustrative to consider that a recent large-scale metagenome assembly of cattle
405 rumen metagenomes generated 69,678 non-redundant genes involved in carbohydrate
406 metabolism⁹, while our gene collection comprised substantially more CAZy-annotated
407 gene clusters ($n = 87,573$), even after collapsing at 50% sequence identity. The
408 increased mappability that we achieved across all 5 vertebrate clades when
409 incorporating our gene catalogue in our functional metagenome profiling pipeline
410 demonstrates how our gene collection will likely aid future vertebrate gut metagenome
411 studies (Figure S15).

412 Our assessment of gene cluster abundances in metagenomes from environment
413 and host-associated biomes illuminates how microbiome functioning and taxonomy is
414 distributed across the free-living to obligate symbiont spectrum. Most notably, nearly all
415 prominent functional groups were enriched in both the environment and host-associated
416 biomes, but the specific gene clusters belonged to different taxonomic groups in each
417 biome (Figure 6). For instance, almost all abundant CAZy families were enriched in both
418 the environment and host biomes, but the environment was dominated by
419 Proteobacteria, while Firmicutes, Bacteroidetes, and Actinobacteria gene clusters
420 comprised most host-enriched CAZy families. This suggests the same coarse-level
421 functional groups are present across the free-living to obligate microbe-vertebrate
422 symbiosis lifestyles, but coarse-level taxonomy strongly differs across this spectrum.
423 This pattern largely remained true when we compared enrichment between the
424 Mammalia and non-mammals, suggesting that taxonomic differences prevail over
425 functional differences in regards to host specificity, at least over broad-scale vertebrate
426 evolutionary distances.

427 In conclusion, our large-scale metagenome assembly of both MAGs and coding
428 sequences from a highly diverse collection of vertebrates greatly expands the known
429 taxonomic and functional diversity of the vertebrate gut microbiome. We have
430 demonstrated that both taxonomic and functional metagenome profiling of the
431 vertebrate gut is improved by our MAG and gene catalogues, which will aid future
432 investigations of the vertebrate gut microbiome. Moreover, our collection can help guide
433 natural product discovery and bioprospecting of novel carbohydrate-active enzymes,
434 along with modeling AMR transmission among reservoirs. By characterizing the
435 distribution of MAGs and microbial genes across environment and host biomes, we
436 gained insight into how taxonomy and function differ along the free-living to obligate
437 symbiosis lifestyle spectrum. We must note that our metagenome assembly dataset is
438 biased toward certain animal clades, which likely impacts these findings. As
439 metagenome assembly becomes more commonplace for studying the vertebrate gut
440 microbiome, bias toward certain vertebrates (*e.g.*, humans) will decrease, and thus
441 allow for a more comprehensive reassessment of our findings.

442 **Acknowledgements**

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

451 We would like to thank the following collaborators for their huge efforts in sample
452 and data collection: Mario Baldi, School of Veterinary Medicine, Universidad Nacional
453 de Costa Rica; Wolfgang Vogl and Frank Radon, Konrad Lorenz Institute of Ethology
454 and Biological Station Illmitz; Endre Sós and Viktor Molnár, Budapest Zoo; Ulrike
455 Streicher, Conservation and Wildlife Management Consultant, Vietnam; Katharina Mahr,
456 Konrad Lorenz Institute of Ethology, University of Veterinary Medicine Vienna and
457 Flinders University Adelaide, South Australia; Peggy Rismiller, Pelican Lagoon
458 Research Centre, Australia; Rob Deaville, Institute of Zoology, Zoological Society of
459 London; Alex Lécu, Muséum National d'Histoire Naturelle and Paris Zoo; Danny
460 Govender and Emily Lane, South African National Parks, Sanparks; Fritz Reimoser,
461 Research Institute of Wildlife Ecology, University of Veterinary Medicine Vienna; Anna
462 Kübber-Heiss and Team, Pathology, Research Institute of Wildlife Ecology, University of
463 Veterinary Medicine Vienna; Nikolaus Eisank, Nationalpark Hohe Tauern, Kärnten; Attila

464 Hettyey and Yoshan Moodley, Konrad Lorenz Institute of Ethology, University of
465 Veterinary Medicine Vienna; Mansour El-Matbouli and Oskar Schachner, Clinical Unit of
466 Fish Medicine, University of Veterinary Medicine; Barbara Richter, Institute of Pathology
467 and Forensic Veterinary Medicine, University of Veterinary Medicine Vienna; Hanna
468 Vielgrader and Zoovet Team, Schönbrunn Zoo; Reinhard Pichler, Herberstein Zoo. We
469 explicitly thank the Freek Venter of South African National Parks and the National
470 Zoological Gardens of South Africa for granting access to their Parks for sample
471 collection.

472 **Author Contributions**

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

477 **Competing Interest Statement**

478 No conflicts of interest declared.

479 **References**

- 480 1. Nayfach, S., Rodriguez-Mueller, B., Garud, N. & Pollard, K. S. An integrated
481 metagenomics pipeline for strain profiling reveals novel patterns of bacterial
482 transmission and biogeography. *Genome Res.* **26**, 1612–1625 (2016).
- 483 2. Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome
484 research. *BMC Biol.* **17**, 48 (2019).
- 485 3. Wang, W.-L. *et al.* Application of metagenomics in the human gut microbiome.
486 *World J. Gastroenterol.* **21**, 803–814 (2015).
- 487 4. Zou, Y. *et al.* 1,520 reference genomes from cultivated human gut bacteria enable
488 functional microbiome analyses. *Nat. Biotechnol.* **37**, 179–185 (2019).
- 489 5. Forster, S. C. *et al.* A human gut bacterial genome and culture collection for
490 improved metagenomic analyses. *Nat. Biotechnol.* **37**, 186–192 (2019).
- 491 6. Mukherjee, S. *et al.* 1,003 reference genomes of bacterial and archaeal isolates
492 expand coverage of the tree of life. *Nat. Biotechnol.* **35**, 676–683 (2017).
- 493 7. Almeida, A. *et al.* A unified sequence catalogue of over 280,000 genomes obtained
494 from the human gut microbiome. *bioRxiv* 762682 (2019) doi:10.1101/762682.
- 495 8. Huang, P. *et al.* The chicken gut metagenome and the modulatory effects of
496 plant-derived benzylisoquinoline alkaloids. *Microbiome* **6**, 211 (2018).
- 497 9. Stewart, R. D. *et al.* Assembly of 913 microbial genomes from metagenomic
498 sequencing of the cow rumen. *Nat. Commun.* **9**, 870 (2018).
- 499 10. Riiser, E. S. *et al.* Switching on the light: using metagenomic shotgun sequencing to
500 characterize the intestinal microbiome of Atlantic cod. *Environ. Microbiol.* **21**,

- 501 2576–2594 (2019).
- 502 11. Gibson, K. M. *et al.* Gut microbiome differences between wild and captive black
503 rhinoceros - implications for rhino health. *Sci. Rep.* **9**, 7570 (2019).
- 504 12. Youngblut, N. D. *et al.* Host diet and evolutionary history explain different aspects of
505 gut microbiome diversity among vertebrate clades. *Nat. Commun.* **10**, 2200 (2019).
- 506 13. Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time.
507 *Nat. Commun.* **9**, 2542 (2018).
- 508 14. Steinegger, M., Mirdita, M. & Söding, J. Protein-level assembly increases protein
509 sequence recovery from metagenomic samples manyfold. *Nat. Methods* **16**,
510 603–606 (2019).
- 511 15. Karasov, T. L. *et al.* Arabidopsis thaliana and Pseudomonas Pathogens Exhibit
512 Stable Associations over Evolutionary Timescales. *Cell Host Microbe* **24**,
513 168–179.e4 (2018).
- 514 16. Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: estimating species
515 abundance in metagenomics data. *PeerJ Comput. Sci.* **3**, e104 (2017).
- 516 17. de la Cuesta-Zuluaga, J., Ley, R. E. & Youngblut, N. D. Struo: a pipeline for building
517 custom databases for common metagenome profilers. *Bioinformatics* (2019)
518 doi:10.1093/bioinformatics/btz899.
- 519 18. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and
520 metatranscriptomes. *Nat. Methods* **15**, 962–968 (2018).
- 521 19. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new
522 method for improved phylogenetic and taxonomic placement of microbes. *Nat.*
523 *Commun.* **4**, 2304 (2013).
- 524 20. Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining
525 pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).
- 526 21. Hannigan, G. D. *et al.* A deep learning genome-mining strategy for biosynthetic
527 gene cluster prediction. *Nucleic Acids Res.* **47**, e110 (2019).
- 528 22. Navarro-Muñoz, J. C. *et al.* A computational framework to explore large-scale
529 biosynthetic diversity. *Nat. Chem. Biol.* **16**, 60–68 (2020).
- 530 23. Breitwieser, F. P., Baker, D. N. & Salzberg, S. L. KrakenUniq: confident and fast
531 metagenomics classification using unique k-mer counts. *Genome Biol.* **19**, 198
532 (2018).
- 533 24. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
534 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 535 25. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to
536 classify genomes with the Genome Taxonomy Database. *Bioinformatics* (2019)
537 doi:10.1093/bioinformatics/btz848.
- 538 26. Allen, H. K. *et al.* Call of the wild: antibiotic resistance genes in natural
539 environments. *Nat. Rev. Microbiol.* **8**, 251–259 (2010).
- 540 27. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
541 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
- 542 28. Sachs, J. L., Skophammer, R. G. & Regus, J. U. Evolutionary transitions in
543 bacterial symbiosis. *Proc. Natl. Acad. Sci. U. S. A.* **108 Suppl 2**, 10800–10807
544 (2011).
- 545 29. Ebert, D. The Epidemiology and Evolution of Symbionts with Mixed-Mode

- 546 Transmission. *Annu. Rev. Ecol. Evol. Syst.* **44**, 623–643 (2013).
- 547 30. Robinson, C. D. *et al.* Experimental bacterial adaptation to the zebrafish gut reveals
548 a primary role for immigration. *PLoS Biol.* **16**, e2006893 (2018).
- 549 31. Adnani, N., Rajski, S. R. & Bugni, T. S. Symbiosis-inspired approaches to antibiotic
550 discovery. *Nat. Prod. Rep.* **34**, 784–814 (2017).
- 551 32. Donia, M. S. *et al.* A systematic analysis of biosynthetic gene clusters in the human
552 microbiome reveals a common family of antibiotics. *Cell* **158**, 1402–1414 (2014).
- 553 33. Donia, M. S. & Fischbach, M. A. Small molecules from the human microbiota.
554 *Science* **349**, 1254766 (2015).
- 555 34. Milshteyn, A., Colosimo, D. A. & Brady, S. F. Accessing Bioactive Natural Products
556 from the Human Microbiome. *Cell Host Microbe* **23**, 725–736 (2018).
- 557 35. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes
558 substantially expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).
- 559 36. Xiao, L. *et al.* A reference gene catalogue of the pig gut microbiome. *Nat Microbiol*
560 **1**, 16161 (2016).
- 561 37. Coelho, L. P. *et al.* Similarity of the dog and human gut microbiomes in gene
562 content and response to diet. *Microbiome* **6**, 72 (2018).
- 563 38. Pan, H. *et al.* A gene catalogue of the Sprague-Dawley rat gut metagenome.
564 *Gigascience* **7**, (2018).