

1 A curated data resource of 214K
2 metagenomes for characterization of the
3 global resistome
4

5 **Short title:** Abundance data of 214K metagenomes of the global resistome
6

7 **Authors:**

8 Hannah-Marie Martiny^{*1},

9 Patrick Munk¹,

10 Christian Brinch¹,

11 Frank M. Aarestrup¹,

12 Thomas N. Petersen¹
13

14 **Affiliations:**

15 ¹Research Group for Genomic Epidemiology, Technical University of Denmark, Kongens
16 Lyngby, Denmark
17

18 * Corresponding author

19 E-mail: hanmar@food.dtu.dk (HMM)
20

21 Abstract

22 The growing threat of antimicrobial resistance (AMR) calls for new epidemiological
23 surveillance methods, as well as a deeper understanding of how antimicrobial resistance
24 genes (ARGs) have transmitted around the world. The large pool of sequencing data
25 available in public repositories provides an excellent resource for monitoring the temporal
26 and spatial dissemination of AMR in different ecological settings. However, only a limited
27 number of research groups globally have the computational resources allowing analyses of
28 such data. We retrieved 442 Tbp of sequencing reads from 214,095 metagenomic samples
29 from the European Nucleotide Archive (ENA) and aligned them using a uniform approach
30 against ARGs and 16S/18S rRNA genes. Here, we present the results of this extensive
31 computational analysis and share the counts of reads aligned. Over $6.76 \cdot 10^8$ read
32 fragments were assigned to ARGs and $3.21 \cdot 10^9$ to rRNA genes, where we observed distinct
33 differences in both the abundance of ARGs and the link between microbiome and resistome
34 compositions across various sampling types. This collection is another step towards
35 establishing a global surveillance of AMR and can serve as a resource for further research
36 into the environmental spread and dynamic changes of ARGs.

38 Introduction

39 The vast amount of genomic data available in public data repositories is a unique and
40 potentially important resource for doing research and genomic surveillance of antimicrobial
41 resistance (AMR). Using these datasets collected from locations all over the world across
42 different years and from various sampling sources might further aid our understanding of
43 the emergence and distribution of antimicrobial resistance genes (ARGs).

44
45 The sharing of genomic sequence data to one of the available repositories is today a major
46 and often mandatory step in peer-reviewed journals, for which several repositories were
47 created by the members of the International Nucleotide Sequence Database Collaboration
48 (INSDC)¹, including the European Nucleotide Archive (ENA)². The number of sequencing data
49 available at ENA continues to increase with an estimated doubling time of 18 months
50 (<https://www.ebi.ac.uk/ena/browser/about/statistics>, accessed 2022-03-08).

51 Several approaches for how to analyze genomic data depending on the sample types are
52 already well established.
53 However, the exploration of these resources are often restricted to a few research groups
54 only, since both sufficient skills in bioinformatics and access to high-performing computer
55 resources are needed to handle the large amount of available data.
56 Existing collections of analyzed datasets tend to focus on either specific sample sources,
57 such as humans^{3,4}, marine⁵, or urban sewage^{6,7}, or focus on specific genera⁸. Especially the
58 COVID19 pandemic has highlighted the value of data sharing to trace the spread and
59 evolution of the virus⁹. Despite the attempts to standardize the analysis workflows of these
60 databases, they are limited in their ability to generalize across environments and locations.
61 A recent study¹⁰ has shared a searchable collection of 661K bacterial genomes for exploring
62 the global bacterial diversity across different origins, providing an easy-to-access resource
63 for genomic research. While this is an impressive data-sharing effort, the authors did not
64 include metagenomic samples in their pipeline. Metagenomic techniques aim to sequence
65 all DNA in a sample and can be used to characterize the microbiome in different
66 environments^{11,12}, discover novel organisms¹³, monitor disease^{14,15} and specific genes, such
67 as ARGs^{5,6,16}.

68
69 Here, we present a large-scale metagenomic analysis of 214,095 metagenomic samples
70 retrieved from ENA. We have carried out an assembly-free approach by aligning sequencing
71 reads against ARGs and 16S/18S ribosomal RNA genes. We have previously published an in-
72 depth analysis of the distribution of mobilized colistin resistance¹⁷ based on those data.
73 Now we both share the entire collection of mapping results and showcase how to
74 characterize the global resistome and microbiome with this dataset. The curated metadata
75 and mapping results are available at <https://doi.org/10.5281/zenodo.6519844> and
76 documentation at <https://hmmartiny.github.io/mARG/Tables.html>.

78 Materials and Methods

79 Retrieval of metagenomes

80 We retrieved metagenomic datasets from the European Nucleotide Archive (ENA)²
81 uploaded between 2010-01-01 and 2020-01-01 that had library source as 'METAGEOMIC'

82 and library strategy of 'WGS'. We collected 214,095 sequencing runs from 146,732 samples
83 from 6,307 projects corresponding to 442 Tbp of raw reads taking up 300 TB of storage. The
84 associated metadata for each sample was also retrieved.

85

86 [Preprocessing and mapping of sequencing reads](#)

87 The retrieved raw FASTQ reads were trimmed and aligned against reference sequences, as
88 outlined in Martiny (2022)¹⁷. In brief, we used FASTQC v.0.11.15

89 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) for read quality checking

90 and BBduk2 v.36.49¹⁸ for trimming the raw sequencing reads. With the k-mer based

91 alignment tool KMA 1.2.21¹⁹, the trimmed reads were mapped against reference sequences

92 from two different databases: The AMR gene database ResFinder20 (downloaded 25-01-

93 2020), which contained 3,085 sequences, and the ribosomal rRNA Silva21 gene database

94 (version 138, downloaded 16-01-2020), which had 2,225,272 template sequences with more

95 than 88% of them being 16/18S rRNA genes. Data retrieval, quality checking, trimming, and

96 read alignments were done using the Danish National Supercomputer for Life Sciences

97 (<https://www.computerome.dk/>).

98

99 [Standardization of metadata](#)

100 The following attributes for each metagenome were standardized: sampling location,

101 sampling host or environment (referred to as a host below), and sampling date.

102

103 To standardize the label for sampling locations, we looked at the values entered in the two

104 fields 'country' and 'location'. First, the latitude and longitude coordinates were mapped to

105 a country using the Python library Shapely 1.7.1²² to find the matching area defined in one

106 of the three public domain map datasets (countries, marine, and lakes) available in the

107 Natural Earth Data collection. If the lookup failed or the coordinates were not given, the

108 second step was to match the text attribute in the country label to ISO 3166 country codes

109 with a fuzzy search with the Python library PyCountry 20.7.3

110 (<https://github.com/flyingcircusio/pycountry>). Finally, if the two lookup searches did not

111 yield a match, we did a manual lookup of the country labels to standardize the text.

112

113 For the standardization of host labels, we mapped the taxonomic id given by the attribute
114 'host_tax_id' to the NCBI Taxonomy database²³, or if the feature was missing, the 'tax_id'
115 was used instead.

116

117 Since the only way to curate entered collection dates is to look up suspicious dates in
118 published studies manually, and that was deemed too time-intensive, we decided to replace
119 dates entered as later than 2020-01-01 in the sample attribute field 'collection_date' with
120 the missing value NULL.

121

122 [Measuring the abundance of ARGs](#)

123 Since we report the fragment count aligned to each reference gene, the mapping results are
124 compositional and should be treated as such²⁴. In the simplest form, the ARG abundance for
125 a sample or sample group can be calculated as the log-ratio of the count of reads, n_i ,
126 aligned to each ARG i over the total sum of rRNA read fragments n_B :

$$127 \quad x = [n_1, n_2, \dots, n_D, n_B] , \quad i = 1..D$$

$$128 \quad \text{Abundance}(x) = \left[\log \frac{n_1}{n_B}, \log \frac{n_2}{n_B}, \dots, \log \frac{n_D}{n_B} \right]$$

129 where D is the number of ARGs and $n_B = \frac{\sum_j^{D_B} n_j}{1 \cdot 10^6}$ with D_B being the number of rRNA genes.

130 Each ARG count n_i has been adjusted with the length of the gene in kilobases.

131

132 The relative abundance resistance classes were calculated as the proportion of ARG
133 resistance assigned to different classes and scaled with $\kappa = 100$:

$$134 \quad \text{Relative abundance}(x) = \frac{\kappa}{\sum n_i} n_i$$

135 [Diversity measurements](#)

136 Besides the read abundance values, we report the species richness, Shannon diversity
137 index²⁵, and the Gini-Simpson²⁶ diversity index of read counts of ARGs, genera, and phyla
138 per sample. Species richness is the number of different genes or taxonomic groups present
139 in the sample with at least one read fragment aligned.

140

141 The Shannon index (H') were calculated using the proportions of reads $p_i = \frac{n_i}{\sum n}$:

$$142 \quad H' = - \sum_{i=1}^R p_i \ln p_i$$

143 whereas the Gini-Simpson index (GS) was calculated using the read counts $n = [n_1, \dots, n_D]$
144 and $N = \sum n$ is the total count of reads for the group:

$$145 \quad GS = 1 - \frac{\sum n_i \cdot (n_i - 1)}{N \cdot (N - 1)}$$

146 Together with these two indices, we also report the sample-wise unique number of
147 templates or taxonomic groups matched.

148

149 [Code and data availability](#)

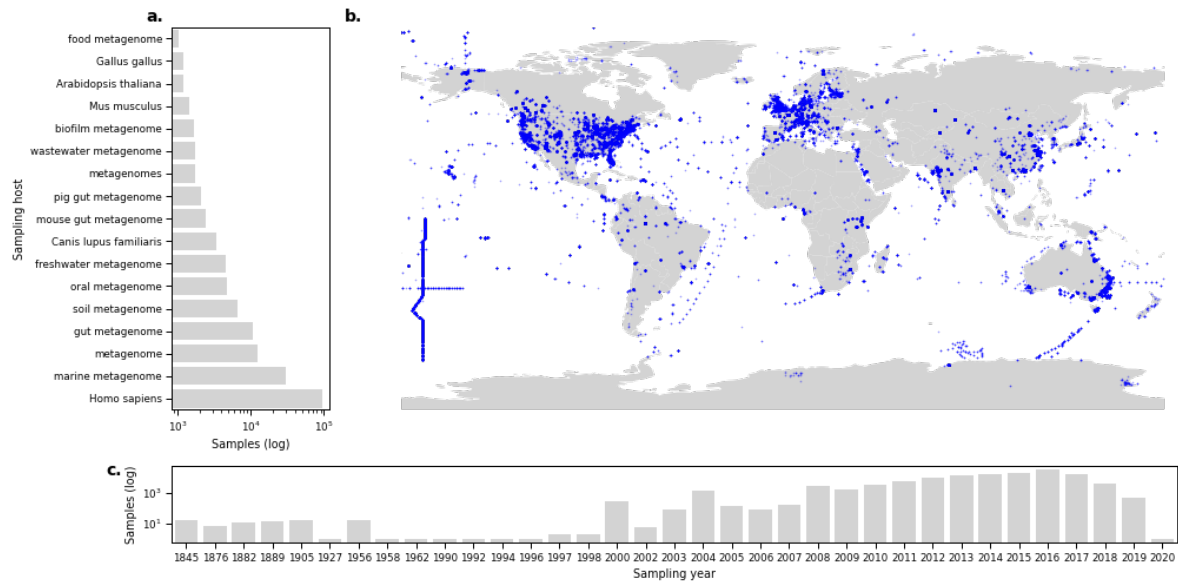
150 The code to produce the figures is available at <https://github.com/hmmartiny/mARG>. The
151 data has been deposited at <https://doi.org/10.5281/zenodo.6519844>, and documentation
152 of the various tables be accessed at <https://hmmartiny.github.io/mARG>.

153 [Results](#)

154 Here, we present a large-scale mapping of 442 Tbp of raw reads of 214,095 metagenomic
155 samples suitable for analyzing the distribution of acquired antimicrobial resistance genes
156 and 16S/18S rRNA genes. Furthermore, we have spent considerable effort standardizing
157 three main sample attributes: sampling date, location, and source. To facilitate easy access
158 and usage, we have shared the mapping results and corrected metadata in three different
159 data formats (TSV, HDF, and MySQL dumps). We also provide tutorials with code examples
160 in R and Python on using the data in different scenarios. Data files are all available at
161 <https://doi.org/10.5281/zenodo.6519844>.

162

163 By collecting the sequencing reads from ENA, we could also verify the inherited bias of
164 specific sample types or sources being overrepresented simply due to the availability in the
165 public repository. While the 214,095 metagenomic datasets were collected from 797
166 different hosts, most were either of human or marine origin (Figure 1a). A similar skewed
167 geographical distribution towards European and North American countries was observed in
168 the sampling locations (Figure 1b). The distribution of samples according to the sampling
169 year reveals that a considerable number were collected between 2010 and 2020 (Figure 1c).



170

171

172

173

174

Figure 1: Distribution of metagenomes reveals the overrepresentation of samples from specific sources.

a. Number of samples grouped per sampling host, where only hosts with more than 1000 samples are plotted. **b.** Sample locations for metagenomes with available GPS coordinates; each marker is a sample. **c.** Year of which a sample was collected.

175

176

177

178

179

180

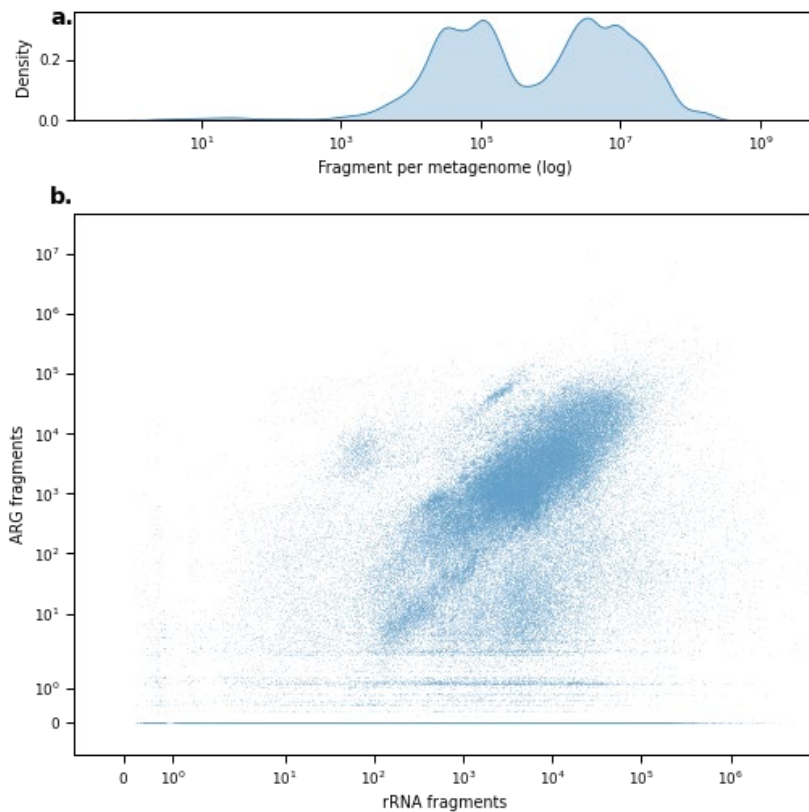
181

182

183

184

Of the more than $1.8 \cdot 10^{12}$ raw sequencing reads, corresponding to 442.1 Tbp, 93% of the reads were generated using Illumina sequencing technologies (Figure S1). We mapped over $1.69 \cdot 10^{12}$ trimmed read fragments, with a median of 784,748 fragments per sample (range 1 – 916,901,400) (Figure 2a). 0.04% of all read fragments could be aligned to ARGs, and 0.19% to rRNA genes. The number of ARG fragments aligned increased with the number of aligned rRNA fragments, although for 34% of the samples, we did not find any ARGs despite having read fragments aligning to 16S rRNA genes (Figure 2b). The microbial differences in the different sampling origins were highlighted in the number of aligned fragments (Figure S3).



185

186

187

Figure 2: Distribution of available and aligned fragments. a. Density distribution of available fragments per sample. b. The distribution compares the number of fragments mapped to rRNA genes and ARGs.

188 The global abundance of antimicrobial resistance

189 To measure the global distribution of ARGs and the composition of the resistome, we
190 calculated the abundance of ARGs as the log-ratio of ARG fragments over summed rRNA
191 sequence fragments. Almost all of the template sequences from the ResFinder database had
192 at least one fragment aligned, and only 94 ARGs had no hits (Figure S2). The median
193 observed resistance load per metagenomic sample was 11.74 (log range: -1.45 to 23.52)
194 (Figure 3a), which appeared to be mainly dependent on the geographic origin and
195 environment (Figure 3b-d) and not on which year the sample was taken. For example,
196 samples originating from locations within Europe showed similar abundance levels for most
197 of the samples but with several outliers, whereas multiple samples from locations in the
198 Oceania region had a much broader load distribution with few outliers (Figure 3c).

199

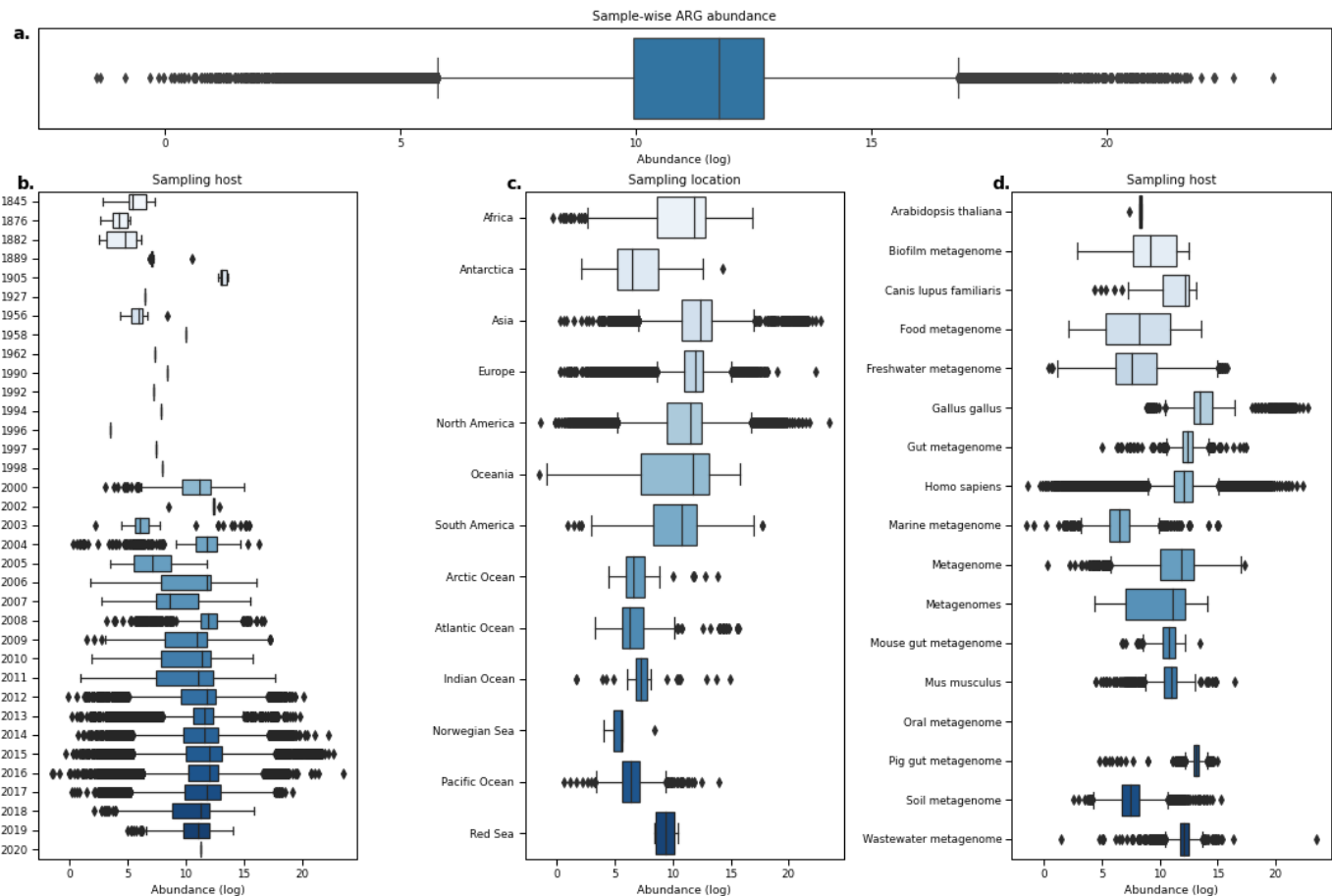


Figure 3: **Boxplots of ARG abundances in metagenomic samples show that levels vary across different origins.**

a. Distribution of ARG abundance per sample. b. Distribution of sample-wise ARG abundance grouped by sampling year. c. Sample-wise ARG abundance per sampling location. d. Sample-wise ARG abundance grouped by hosts. Only hosts with more than 1000 metagenomes analyzed are shown.

200
201
202
203
204

205 While the distribution of sample-wise resistance loads illustrates the high variability in this
206 data collection (Figure 3), we saw that once we stratified the relative ARG read proportions
207 per resistance class and sample type, there were clear separations between different groups
208 (Figure 4). For the sampling years with a considerable number of samples available (2004-
209 2019), the relative proportion of classes was relatively consistent, with Tetracycline reads
210 being the most common, except for a spike of Beta-lactam reads in 2017 (Figure 4a). When
211 looking at the geographic regions, we observed that reads stemming from samples collected
212 from large water bodies had more reads aligned to Aminoglycosides and Beta-lactam classes
213 than land regions, which had more diverse class distributions (Figure 4b), possibly due to
214 that samples from land regions had higher resistance loads overall (Figure 3a). Once we
215 stratified by sampling host or source, the distribution of resistance classes was very
216 dependent on the group, as seen by the high proportion of read fragments aligned to, for

217 example, Phenicol for marine and soil samples and Tetracycline reads being highly prevalent
 218 in mice (*Mus musculus*) samples (Figure 4c).

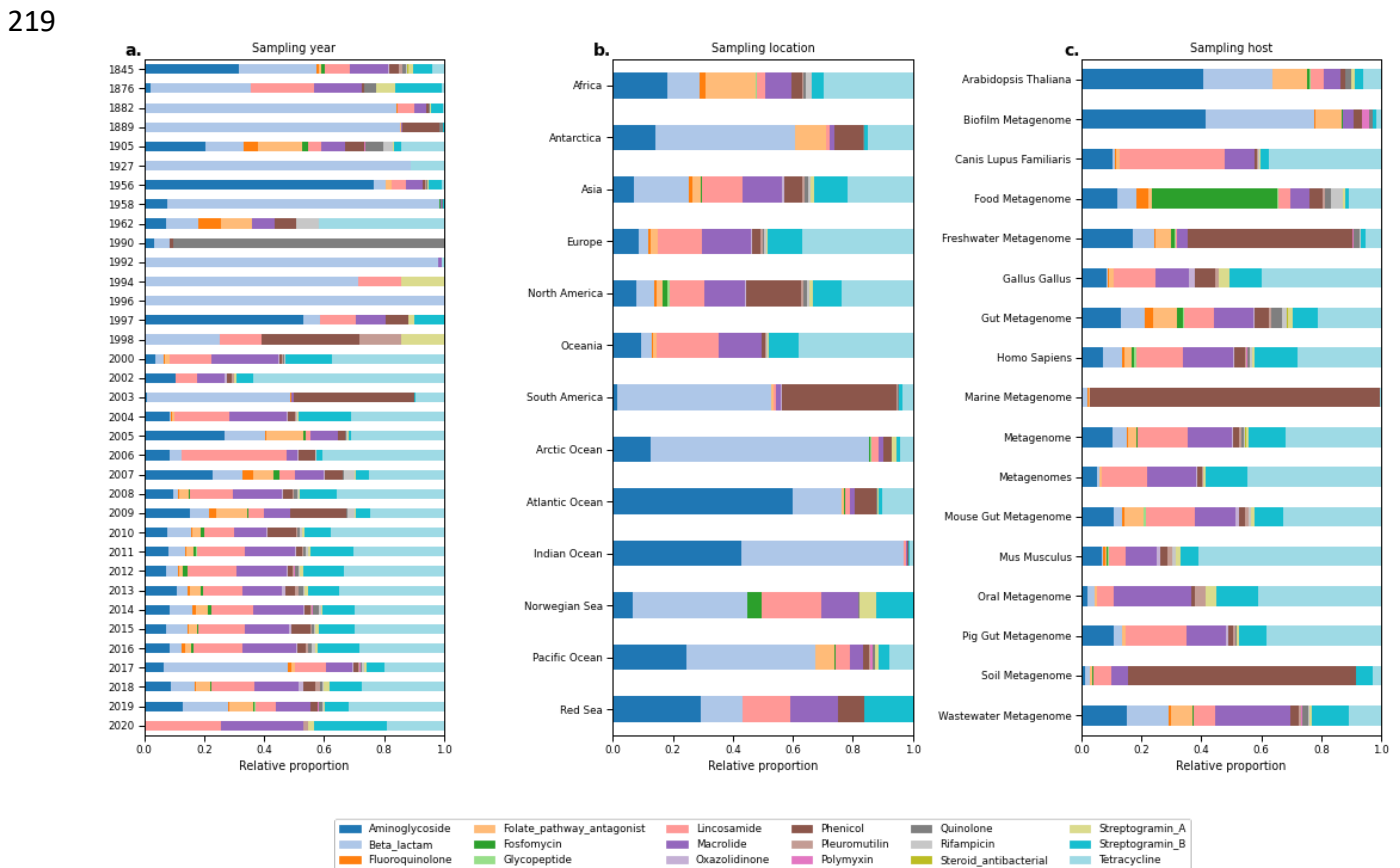
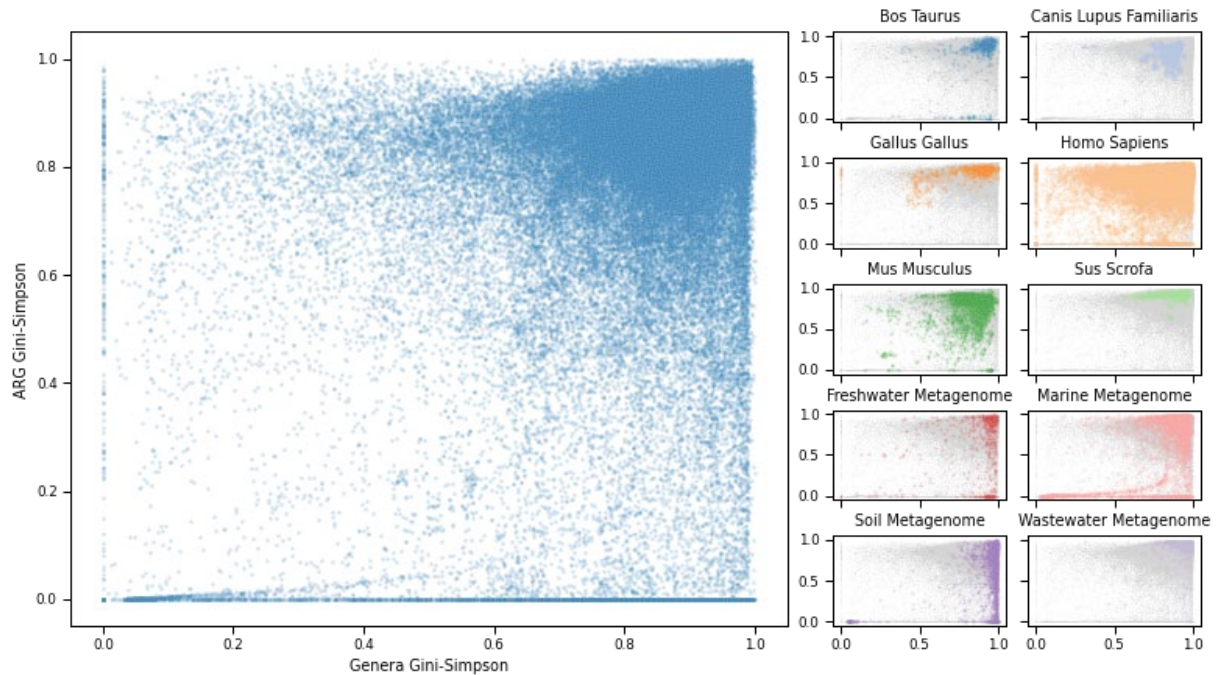


Figure 4: Composition of reads assigned to ARGs from different resistance classes grouped by sampling origin. a. Grouped by sampling year. b. Grouped per sampling location. c. Grouped per sampling host. Only hosts with more than 1000 metagenomes analyzed are shown.

224 Linking the microbiome diversity with resistance diversity

225 The relationship between the diversity of the microbiome and the resistance genes was
 226 quantified by calculating the species richness and two alpha diversity measurements
 227 (Shannon and Gini-Simpson) on ARG levels and phyla and genera taxonomic levels. We saw
 228 a general trend of increased diversity of the microbiome also meant an increased ARG
 229 diversity (Figure 5, Figure S4). Although, the relationship between genera and ARG diversity
 230 indexes further characterized by sampling sources revealed a higher differentiation,
 231 suggesting that increased diversity of microbes in, for example, soil samples does not
 232 necessarily lead to a higher diversity of resistance genes. Contrarily, the chicken (*Gallus*
 233 *gallus*) samples showed that they still had elevated ARG diversity despite having lower
 234 microbial diversity (Figure 5).

235



236
237
238
239

Figure 5: **The genus-ARG diversity relationship for all metagenomic samples.** The Gini-Simpson diversity indexes were calculated on genus categories (*x*-axis) compared to ARG levels (*y*-axis). Left: Scatterplot of all samples. Right: Samples colored by selected host or environmental origins.

240 Discussion

241 Global surveillance of AMR based on genomics continues to become more accessible due to
242 the advancement in NGS technologies and the practice of sharing raw sequencing data in
243 public repositories. Standardized pipelines and databases are needed to utilize these large
244 data volumes for tracking the dissemination of AMR. We have uniformly processed the
245 sequencing reads of 214,095 metagenomes for the abundance analysis of ARGs.
246 Our data sharing efforts enable users to perform abundance analyses of individual ARGs, the
247 resistome, and the microbiome across different environments, geographic locations, and
248 sampling years.

249

250 We have given a brief characterization of the distribution of ARGs according to the
251 collection of metagenomes. However, in-depth analyses remain to be performed to
252 investigate the influence of temporal, geographical, and environmental origins on the
253 dissemination and evolution of antimicrobial resistance. For example, analyzing the spread
254 of specific ARGs across locations and different environments could reveal new transmission
255 routes of resistance and guide the design of intervention strategies to stop the spread.
256 Another use of the data collection could be to explore how the changes in microbial
257 abundances affect and are affected by the resistome. Furthermore, our coverage statistics

258 of reads aligned to ARGs could be used to investigate the rate of new variants occurring in
259 different reservoirs. Even though we have focused on the threat of antimicrobial resistance,
260 potential applications of this resource can be to look at the effects of, e.g., climate changes
261 on microbial compositions.

262

263 We recommend that potential users consider all the confounders present in this data
264 collection in their statistical tests and modeling workflows, emphasizing that the
265 experimental methods and sequencing platforms dictate the obtained sequencing reads.
266 Furthermore, it is important to consider the compositional nature of microbiomes²⁷. The
267 reads do not depend on the distribution of genetic material in the sample but the capacity
268 of the sequencing platform^{24,28}. Various statistical methods already exist that consider the
269 compositionality^{24,29,30}.

270

271 With this data resource, we have taken a step towards enabling the scientific community to
272 utilize the wealth of information in these metagenomic samples to broaden our
273 understanding of the dissemination of antimicrobial resistance and changes in microbiomes
274 at both local and global scales through time and environments.

275 References

- 276 1. Arita, M., Karsch-Mizrachi, I. & Cochrane, G. The international nucleotide sequence
277 database collaboration. *Nucleic Acids Res.* **49**, D121 (2021).
- 278 2. Leinonen, R. *et al.* The European nucleotide archive. *Nucleic Acids Res.* **39**, 44–47
279 (2011).
- 280 3. Shao, L., Liao, J., Qian, J., Chen, W. & Fan, X. MetaGeneBank: a standardized database
281 to study deep sequenced metagenomic data from human fecal specimen. *BMC*
282 *Microbiol.* **21**, 1–12 (2021).
- 283 4. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human
284 gut microbiome. *Nat. Biotechnol.* **39**, 105–114 (2021).
- 285 5. Cuadrat, R. R. C., Sorokina, M., Andrade, B. G., Goris, T. & Dávila, A. M. R. Global
286 ocean resistome revealed: Exploring antibiotic resistance gene abundance and
287 distribution in TARA Oceans samples. *Gigascience* **9**, 1–12 (2020).
- 288 6. Hendriksen, R. S. *et al.* Global monitoring of antimicrobial resistance based on

- 289 metagenomics analyses of urban sewage. *Nat. Commun.* (2019). doi:10.1038/s41467-
290 019-08853-3
- 291 7. Fresia, P. *et al.* Urban metagenomics uncover antibiotic resistance reservoirs in
292 coastal beach and sewage waters. *Microbiome* **7**, 1–9 (2019).
- 293 8. Zhou, Z., Alikhan, N. F., Mohamed, K., Fan, Y. & Achtman, M. The Enterobase user’s
294 guide, with case studies on Salmonella transmissions, Yersinia pestis phylogeny, and
295 Escherichia core genomic diversity. *Genome Res.* **30**, 138–152 (2020).
- 296 9. Khare, S. *et al.* GISAID’s Role in Pandemic Response. *China CDC Wkly.* **3**, 1049–1051
297 (2021).
- 298 10. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable
299 snapshot of archived DNA sequences. *PLOS Biol.* **19**, e3001421 (2021).
- 300 11. Fierer, N. *et al.* Cross-biome metagenomic analyses of soil microbial communities and
301 their functional attributes. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 21390–21395 (2012).
- 302 12. Gill, S. R. *et al.* Metagenomic analysis of the human distal gut microbiome. *Science*
303 *(80-)*. **312**, 1355–1359 (2006).
- 304 13. Al-Shayeb, B. *et al.* Clades of huge phages from across Earth’s ecosystems. *Nature*
305 **578**, 425–431 (2020).
- 306 14. Nieuwenhuijse, D. F. *et al.* Setting a baseline for global urban virome surveillance in
307 sewage. *Sci. Rep.* **10**, 1–13 (2020).
- 308 15. Liu, P., Chen, W. & Chen, J. P. Viral Metagenomics Revealed Sendai Virus and
309 Coronavirus Infection of Malayan Pangolins (*Manis javanica*). *Viruses* 2019, Vol. 11,
310 Page 979 **11**, 979 (2019).
- 311 16. Forsberg, K. J. *et al.* Bacterial phylogeny structures soil resistomes across habitats.
312 *Nature* **509**, 612–616 (2014).
- 313 17. Martiny, H.-M. *et al.* Global distribution of mcr gene variants in 214,095 metagenomic
314 samples. *mSystems* (2022). doi:10.1128/msystems.00105-22
- 315 18. Bushnell, B. BBMap. (2014).
- 316 19. Clausen, P. T. L. C., Aarestrup, F. M. & Lund, O. Rapid and precise alignment of raw
317 reads against redundant databases with KMA. *BMC Bioinformatics* **19**, 1–8 (2018).
- 318 20. Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *J.*
319 *Antimicrob. Chemother.* **67**, 2640–2644 (2012).
- 320 21. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: Improved data

- 321 processing and web-based tools. *Nucleic Acids Res.* **41**, 590–596 (2013).
- 322 22. Gillies, S. & Others, A. Shapely: manipulation and analysis of geometric objects.
323 (2007).
- 324 23. Federhen, S. The NCBI Taxonomy database. *Nucleic Acids Res.* **40**, D136–D143 (2012).
- 325 24. Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome
326 datasets are compositional: And this is not optional. *Front. Microbiol.* **8**, 1–6 (2017).
- 327 25. Shannon, C. E. A mathematical theory of communication. *Bell Syst. Tech. J.* **27**, 379–
328 423 (1948).
- 329 26. Jost, L. Entropy and diversity. *Oikos* **113**, 363–375 (2006).
- 330 27. Aitchison, J. The Statistical Analysis of Compositional Data. *J. R. Stat. Soc. Ser. B* **44**,
331 139–160 (1982).
- 332 28. Quinn, T. P. *et al.* A field guide for the compositional analysis of any-omics data.
333 *Gigascience* **8**, 1–14 (2019).
- 334 29. Fernandes, A. D., Macklaim, J. M., Linn, T. G., Reid, G. & Gloor, G. B. ANOVA-Like
335 Differential Expression (ALDEx) Analysis for Mixed Population RNA-Seq. *PLoS One* **8**,
336 (2013).
- 337 30. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data.
338 *PLoS Comput. Biol.* **8**, 1–11 (2012).
- 339

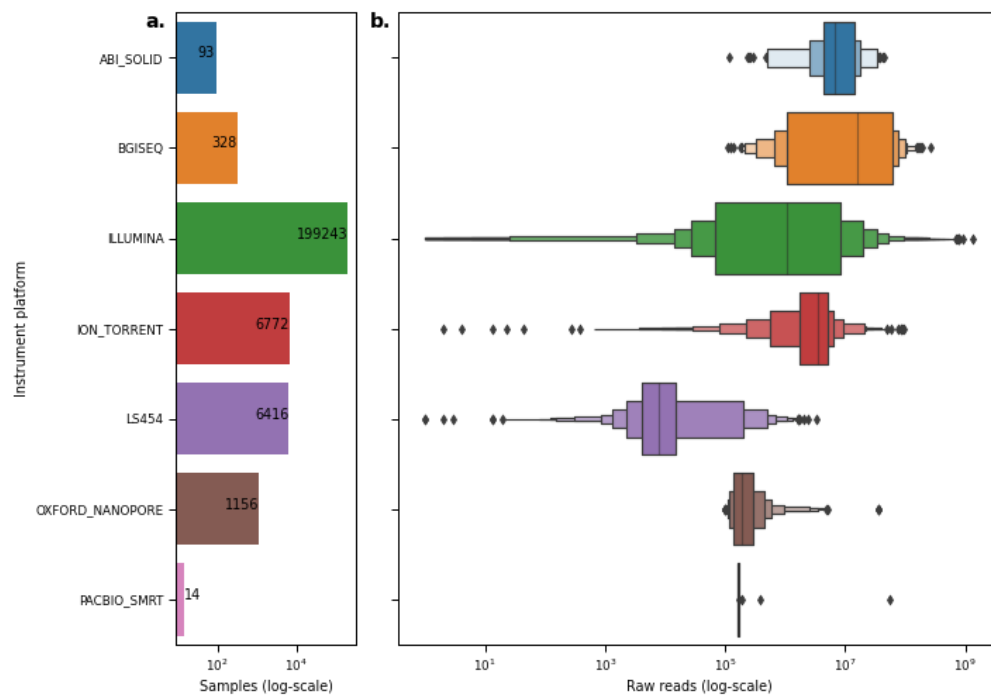


Figure S1: **Distribution of samples per sequencing instrument platform.** **a.** Sample count per platform. **b.** Distribution of raw sequencing read counts per platform.

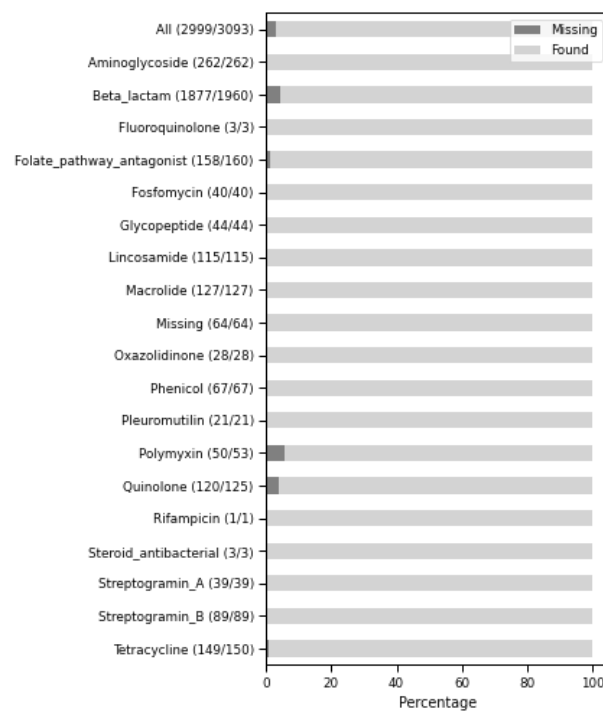


Figure S2: **More than 96% of ARG templates had at least one aligned fragment.** The bars illustrate the percentage of ARGs per resistance class without and with at least one aligned fragment. The parenthesis after each class label contains the number of genes found out of the total available templates.

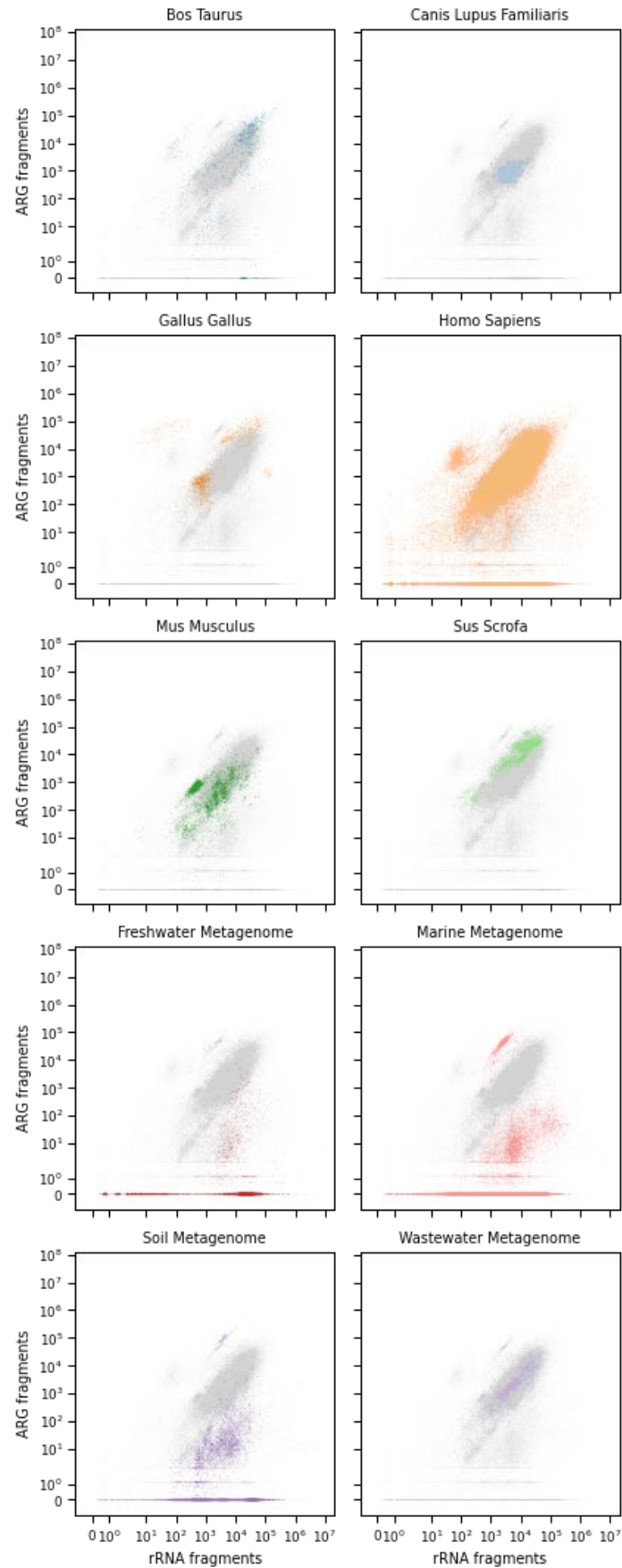


Figure S3: The sample-wise distribution of aligned rRNA fragments and ARG fragments, colored by selected host and environmental sources.

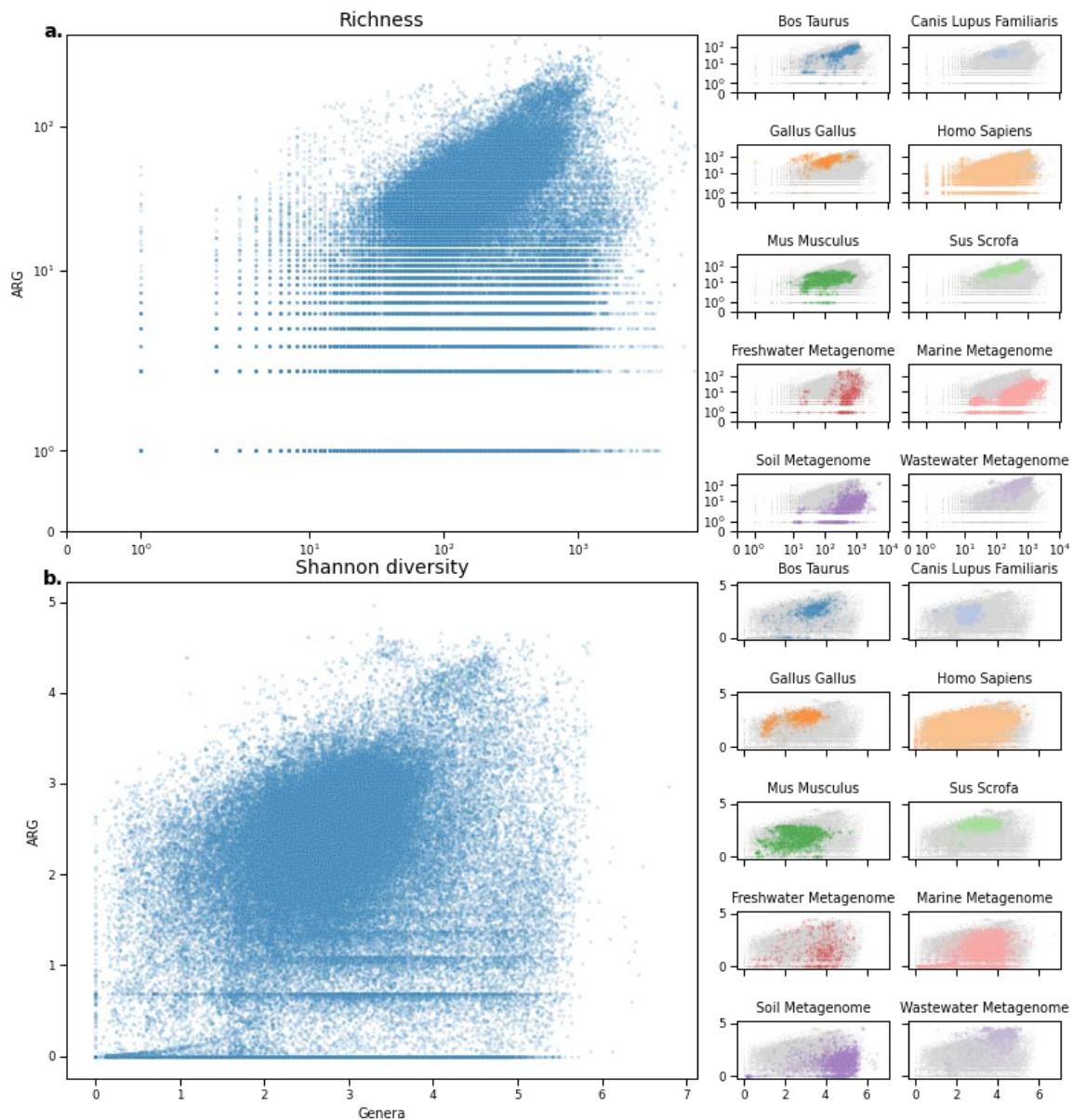


Figure S4: **Additional distributions showing the relationship between ARGs and genera for all metagenomic samples. a.** The richness of genus groups (x-axis) vs. ARG richness (y-axis). **b.** The relationship between Shannon diversity index calculated on genus level (x-axis) and ARGs(y-axis). Right: Samples colored by selected host or environmental origins.