

1 **The distribution of antimicrobial resistance genes across phylogroup, host species and**
2 **geography in 16,000 publicly-available *E. coli* genomes**

3

4 Elizabeth Pursey^{1,2*}, Tatiana Dimitriu¹, William H. Gaze³, Edze R. Westra¹ and Stineke van
5 Houte¹

6

7 ¹Environment and Sustainability Institute, Biosciences, University of Exeter, Penryn,
8 Cornwall, United Kingdom

9 ²Department of Experimental Medicine, University of Lund, Sweden

10 ³European Centre for Environment and Human Health, University of Exeter Medical School,
11 University of Exeter, Penryn, Cornwall, United Kingdom

12

13 *Correspondence to: ellie.pursey@gmail.com

14 Key words: antimicrobial resistance, *E. coli*, phylogroups, hosts, regions

15

16 **Abstract**

17

18 *E. coli* is a highly diverse bacterial species that generates a huge global burden of
19 antimicrobial-resistant infections. A wealth of whole genome sequence data is available on
20 public databases for this species, presenting new opportunities to analyse the distribution of
21 antimicrobial resistance (AMR) genes across its genetic and ecological diversity. We
22 extracted and categorised metadata on host species and geographic location and combined
23 this with *in silico* phylogrouping to describe the characteristics of ~16,000 assembled *E. coli*
24 genomes from the NCBI RefSeq database. We estimated AMR carriage using various
25 metrics: counts of overall genes, multidrug- and extensively drug-resistant categories, and
26 selected β -lactamases of current global concern - *bla*_{CTX-M} and carbapenemase genes. We
27 present estimates of AMR carriage for these metrics by species type (human,
28 agricultural/domestic animal, wild birds and other wild animals), geographic subregion, and
29 across phylogroups. In addition, we describe the distribution of phylogroups within host types
30 and geographic subregions. Our findings show high AMR carriage in commensal-associated
31 phylogroups, agricultural and wild animal hosts and in many subregions. However, we also
32 quantify large biases in sequencing data, the substantial gaps in our knowledge of AMR in
33 many hosts, regions and environmental settings, and the need for systematic sampling to gain
34 a more accurate picture.

35 **Introduction**

36

37 *Escherichia coli* is a remarkably diverse species both genetically and ecologically. It occupies
38 various niches, ranging from a commensal of many warm-blooded organisms, to a globally
39 devastating pathogen, to a free-living environmental bacterium [1]. It is a leading cause of
40 mortality associated with drug-resistant infections and was the pathogen responsible for the
41 most deaths attributed to antimicrobial resistance (AMR) in 2019 [2]. In addition,
42 carbapenem-resistant and extended-spectrum β -lactamase (ESBL)-producing *E. coli* are
43 World Health Organization priority pathogens, for which development of new antibiotics is
44 urgently needed [3]. In light of the ongoing AMR crisis, it is necessary to understand how
45 resistance genes are distributed across the ecological and genetic diversity of this species,
46 highlighting the potential risks they pose across contexts.

47

48 AMR genes (ARGs) are part of the substantial accessory gene content found in *E. coli* [4],
49 and therefore are not distributed uniformly across the species phylogeny. Acquired ARGs are
50 particularly important due to their extensive ability to transmit horizontally in populations on
51 mobile elements such as plasmids, leading to AMR spread across more distant locations and
52 genetic lineages [5]. Some of these ARGs present significant clinical risks in humans,
53 particularly those with current or potential spread across bacterial host genetic backgrounds,
54 including to other species, and along the commensal-pathogen continuum [6].

55

56 Most research has focused on AMR in human-associated and pathogenic isolates [7]. A study
57 focusing on ~1,000 isolates selected to represent the diversity of *E. coli*, mostly including
58 isolates from non-clinical origins, found ARGs were associated with strains from humans and
59 domesticated animals, independent of phylogroup [8]. Therapeutic and growth promoter
60 usage of antibiotics in agriculture is well-documented, and has been linked to high rates of
61 AMR in livestock-associated isolates [9]. Resistance may also spill over into wild animal
62 populations, and there is increasing evidence of wildlife as potential reservoirs of AMR [10].
63 A One Health approach to AMR emphasises the interconnectedness of human, animal and
64 environmental microbiomes [11], with the potential for bidirectional flow of ARGs between
65 animal (wild or domesticated), crop, natural and built environment, and human microbiomes.
66 There are also geographic and socioeconomic biases in studies of AMR in *E. coli*. Whilst
67 antibiotics are used most in high-resource settings, the burden of AMR is estimated to

68 disproportionately affect lower and middle income countries, despite the scarcity of data in
69 these regions [2].

70

71 A wealth of previous work on *E. coli* has investigated AMR trends in specific settings, such
72 as individual hospitals or farms. As sequencing technologies become more accessible and
73 public genomic databases grow, much larger genomics-based studies of AMR in highly-
74 sampled species like *E. coli* are becoming increasingly feasible. For example, a recent large-
75 scale study of over 70,000 *E. coli* genomes found that the number of distinct ARGs varied
76 between phylogroups, and that specific resistance genes were found more frequently in some,
77 though the pattern of resistance to antibiotic classes remained stable across groups [12].

78 Another study curated a collection of ~10,000 *E. coli* isolates sampled from human sources,
79 and identified lineages where >50% of isolates were MDR [13].

80

81 In this work we expand on previous large-scale genomic studies of ARGs in *E. coli* by
82 investigating diversity across different host categories and geographic subregions, as well as
83 phylogroups. We characterise a final dataset of ~16,000 publicly-available assembled *E. coli*
84 genomes, documenting the distribution of ARGs according to several metrics: counts of
85 ARGs, multidrug- and extensively drug-resistant (MDR and XDR) classifications, and
86 presence of clinically-important β -lactamases (*bla*_{CTX-M} and carbapenemases). First, we used
87 model comparison to assess the contributions of host, subregion and phylogroup towards
88 explaining ARG variation in *E. coli*. Then, we used generalised linear models (GLMs) to
89 describe the probabilities of MDR, *bla*_{CTX-M} and carbapenemase presence according to the
90 same predictors. To do so, we used subsampling and resampling approaches to correct for
91 sample size disparities between groupings, and investigate the sensitivity of trends according
92 to subsample sizes.

93

94 **Methods**

95

96 *Genomes*

97

98 The initial complete dataset of 26,881 *E. coli* genomes was retrieved from the National
99 Center for Biotechnology Information (NCBI) RefSeq database in February 2022 in
100 nucleotide fasta format using ncbi-genome-download v0.3.0 ([https://github.com/kblin/ncbi-
101 genome-download](https://github.com/kblin/ncbi-genome-download)).

102

103 *Detection of ARGs, phylogrouping and retrieval of metadata*

104

105 NCBI AMRFinderPlus v3.10.21 was used to predict the identity of acquired ARGs [14].

106 Clermont phylogrouping [15] was performed *in silico* using the EzClermont command-line

107 tool [16]. Snakemake v6.2.1 [17] was used to generate a pipeline for phylogrouping and

108 detecting ARGs, which was run using the University of Exeter's Advanced Research

109 Computing facilities. We removed efflux (*acrF*, *emrD* and *mdtM*) and *blaEC* β -lactamase

110 genes from the dataset, as they were either very common or ubiquitous. As we could not use

111 phenotypic resistance categories, we used the classes included in NCBI AMRFinderPlus

112 output as a proxy for categorising genomes as multidrug-resistant (MDR; ARGs conferring

113 resistance to $\geq 3/20$ antibiotic classes) and extensively drug-resistant (XDR; ARGs conferring

114 resistance to $\geq 10/20$ antibiotic classes) [18].

115

116 Host and location metadata were retrieved and categorised using the Bio.Entrez utilities from

117 Biopython v1.77. We filtered genomes for those with complete metadata for host species and

118 geographic location, as well as those not typed as *E. coli* by phylogrouping, which excluded

119 10,609 genomes, resulting in a final dataset of 16,272 genomes. Geographic locations were

120 split into 20 subregions (Fig. S1) according to Natural Earth data

121 (<https://www.naturalearthdata.com/>).

122

123 All genomes were sorted into the following host species categories: 'Human',

124 'Agricultural/Domestic animals', 'Wild birds' and 'Wild animals'. This was achieved using

125 regular expressions constructed by manually reviewing text in the 'host' field of the

126 biosample data for each accession number. Wild birds were separated from other wild

127 animals due to their comparatively large numbers of genomes, high mobility and potential for

128 ARG spread [19]. Any animal that is most likely to be associated with humans, such as

129 'mouse' or 'canine', was assigned to the 'Agricultural/Domestic animals' category, which

130 also included farm animals such as cows and chickens. Broader responses that could not be

131 reliably classified into any specific group, such as 'Animal', and 'Avian', were not included

132 in further analyses, as well as animals likely to be human food sources such as mussels.

133 Though a conservative approach was taken when assigning genomes to wild animal or bird

134 host groups, these classifications determine the likely category of the host based on species,

135 and cannot rule out cases such as zoo animals.

136

137 *Statistical analysis*

138

139 Data tidying and statistical analysis was done using R v4.0.4, with dplyr v1.0.1, ggplot2
140 v3.3.6 [20] and MetBrewer [21] for colour palettes. Maps additionally used the R packages sf
141 1.0-7 [22] and rnaturalearth v0.1.0.

142

143 In all cases, binomial GLMs were fitted using glmmTMB v1.1.5 [23]. To assess model fit,
144 simulated residuals were plotted and their dispersion tested using DHARMA v0.4.5 [24].
145 Prior to all statistical modelling, subgroupings with sample sizes ≤ 100 (e.g. Central Asia)
146 were removed from the dataset. Sample sizes for each predictor in the full dataset are shown
147 in Figure 1A-C.

148

149 Initially, maximal models were fitted using the entire dataset with MDR, XDR or *bla*_{CTX-M}
150 genes as a binomial response variable (presence/absence) and all explanatory variables
151 included (host, phylogroup and subregion). For carbapenemase genes, the same response
152 variables were used, but only host and phylogroup were included as predictors due to lack of
153 data coverage for many subregions for these comparatively rare genes. Sample sizes for all
154 combinations of response variable and predictor are shown in Tables S1-3. Next, Akaike
155 information criterion (AIC) scores were used to assess the fit of all possible models using the
156 *dredge* function from MuMIn v1.46.0. In all cases, large Δ AIC values (>50) were seen
157 between the maximal model and the next best model (Table S4), emphasising the
158 contribution of all variables to the distribution of AMR gene measures. Model coefficients for
159 all maximal models are shown in Table S5.

160

161 Next, subsampling approaches were used to circumvent disparities in sampling sizes between
162 groups. Separate models were made for each predictor due to the difficulty in 1) taking single
163 representative subsamples spanning all hosts, phylogroups and subregions and 2) estimating
164 means from models with multiple categorical explanatory variables. A subsample was taken
165 for each predictor, giving datasets with equal numbers of genomes from each phylogroup,
166 host or subregion. Several subsample sizes were used to determine the sensitivity of trends to
167 the sample size. Firstly, the dataset was subsampled without replacement down to the lowest
168 group size for each model (e.g. all host categories were sampled to $n=234$, the sample size for
169 wild birds). In addition, samples with replacement were taken with increasing sizes ($n=500$,

170 1000, 2000 and 4000) for hosts, phylogroups and subregions, to investigate the robustness of
171 trends to changing sample sizes. All single predictor models were also run with the full, non-
172 subsampled dataset.

173

174 Model estimated means and 95% confidence intervals (CIs) were calculated using *ggpredict*
175 from the *ggeffects* package [25]. Model coefficients are presented in Table S6 for host, Table
176 S7 for subregion and Table S8 for phylogroup. We did not perform post-hoc comparisons or
177 report *P*-values in the manuscript due to the large number of pairwise comparisons that would
178 be necessary (e.g. 91 comparisons for each model of subregion data). In addition, we did not
179 have *a priori* hypotheses about group differences to test, but rather aimed to describe the
180 patterns in the data with our modelling approaches. Finally, many groups had large sample
181 sizes, which can produce a significant result even for a biologically insignificant effect size
182 [26, 27].

183

184 **Results**

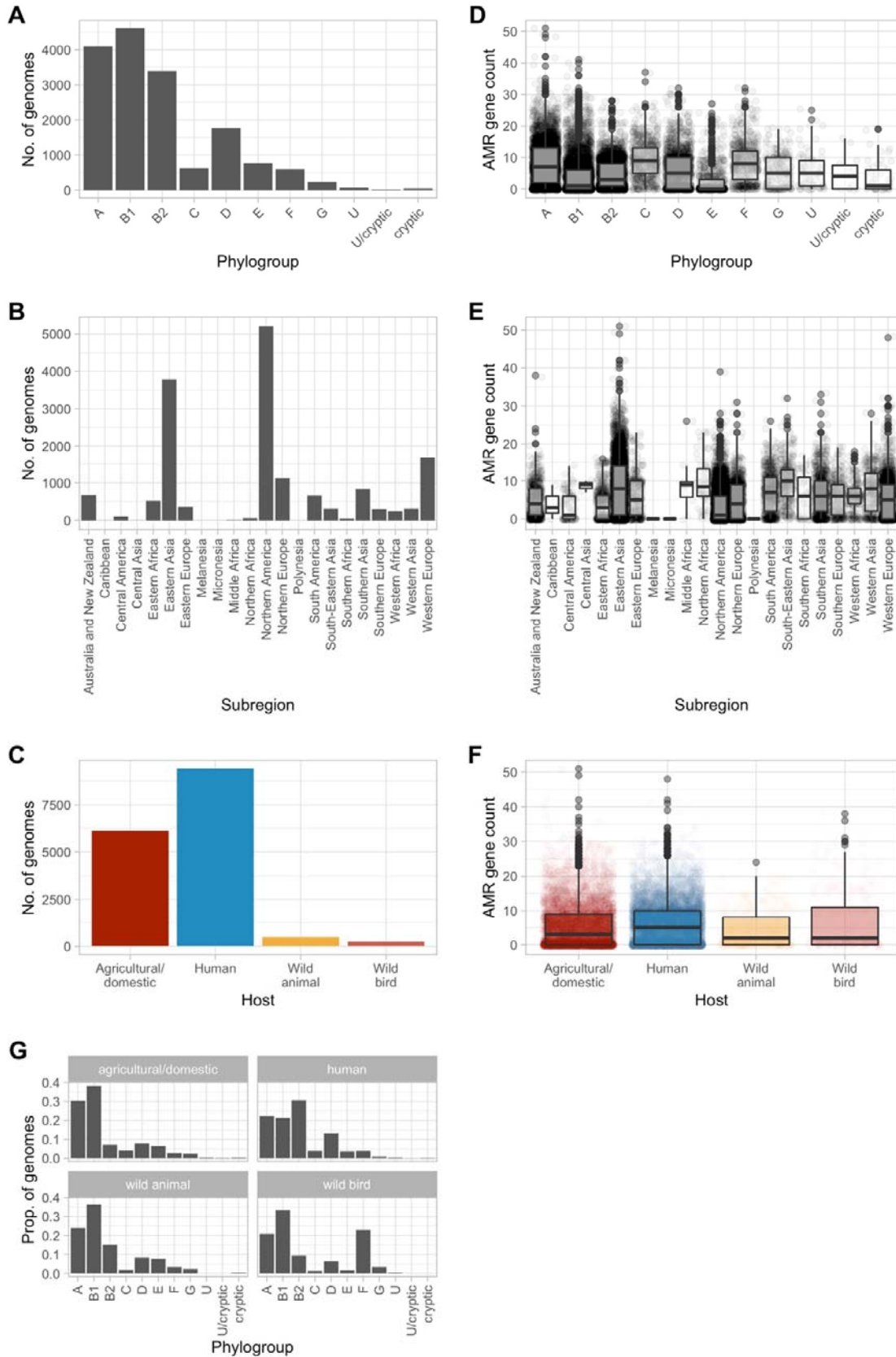
185

186 *Quantification of sample sizes across hosts, phylogroups and subregions in the RefSeq* 187 *dataset*

188

189 Initially, we characterised the RefSeq dataset by looking at sample sizes across phylogroups,
190 host categories, and geographic subregions (Fig. 1A-C, Tables S1-3). Phylogroups B1, A and
191 B2 collectively made up 74% (12,109/16,272) of genomes, with smaller numbers found
192 across C, E, F, G and U (Fig. 1A). Sixty-five genomes were typed as U/cryptic or cryptic.
193 The majority of genomes had a sampling location in North America (n=5,202) or Eastern
194 Asia (n=3,790), whilst many regions were represented poorly such as Micronesia, Melanesia
195 and Polynesia (n=1 for all), Central Asia and Caribbean (n=3 for both) and Middle Africa
196 (n=10) (Fig. 1B). Finally, 96% of the *E. coli* genomes were isolated from humans and
197 agricultural or domestic animals (15,549/16,272) with only 234 from wild birds and 489 from
198 other wild animals (Fig. 1C).

199



200

201 Figure 1 – Total number of genomes per A) phylogroup, B) geographic subregion and C) host species
202 category. Counts of ARGs for D) phylogroup, E) geographic subregion and F) host species category,
203 with raw data plotted below boxplots. G) Proportion of genomes in each phylogroup per host
204 category.

205

206 *ARG counts across the diversity of E. coli sequences*

207

208 We measured counts of ARGs detected across the diversity of genomes to get an overall
209 measure of how AMR carriage was distributed across the dataset (Fig. 1D-F). Overall, the
210 mean number of ARGs per genome for the entire dataset was 5.75. Median ARG counts were
211 highest in phylogroups C, F and A (9, 8 and 7 respectively; Fig. 1D). Meanwhile the lowest
212 median counts were in groups E (0) and B1 (1). A wide range of ARG counts was also seen
213 between and within subregions (Fig. 1E). Among the most well-sampled regions, genomes
214 from Eastern Asia had the highest median ARG count (8), though relatively high numbers
215 were also seen in Western Europe (median ARG count = 5). Additionally, higher median
216 ARG numbers were seen in less represented regions such as Northern (8.5), Middle (9) and
217 Southern Africa (6) and Central Asia (9). High within-group diversity in ARG counts was
218 again seen for host categories, with highest median counts in genomes from humans (5) and
219 agricultural animals (3).

220

221 *Hosts and subregions show different phylogroup patterns*

222

223 We investigated the distribution of *E. coli* phylogroups within hosts (Fig. 1G) and geographic
224 subregions (Fig. S2). *E. coli* isolates from agricultural and domestic animals mostly belonged
225 to phylogroups A and B1 (n=1,863 and n=2,354, respectively). Those from humans were
226 most likely to belong to phylogroup B2 (n=2,870), as well as A and B1 (n=2,071 and
227 n=2,000 respectively). There were also comparatively high numbers of phylogroup D isolates
228 in genomes from human hosts (n=1,237). A higher proportion of phylogroup F isolates were
229 seen in wild birds (n=54), whilst wild animals had slightly higher proportions of B2 isolates
230 (n=74) than agricultural or domesticated animals.

231

232 The frequencies of different phylogroups varied between geographic subregions (Fig. S2).
233 While A, B1 and B2 were typically the largest groups, their relative proportions varied. For
234 example, B1 was more common in Eastern (n=195) and Western (n=94) Africa, B2 was

235 highest in Australia and New Zealand (n=216) and Eastern Europe (n=112), while A was
236 most common in South-Eastern (n=131) and Western (n=117) Asia. Isolates belonging to
237 phylogroup D were also relatively common, with notably higher proportions in regions
238 including Northern Africa (n=17) and Eastern Europe (n=86).

239

240 *Specific measures of AMR carriage*

241

242 Following characterisation of the dataset, we began analysing more specific measures of
243 AMR carriage. First, we categorised genomes as MDR (≥ 1 ARGs conferring resistance to \geq
244 3 classes) and XDR (≥ 1 ARGs conferring resistance to ≥ 10 classes) following detection of
245 402 unique resistance genes conferring predicted resistance to 19 classes. We also looked at
246 the presence of selected β -lactamase genes of current global health concern. We detected 40
247 unique *bla*_{CTX-M} genes across the final dataset, most commonly *bla*_{CTX-M-15} (n=2,536), *bla*_{CTX-}
248 *M-14* (n=1,063) and *bla*_{CTX-M-55} (n=903). Meanwhile, a total of 46 different carbapenemases
249 from the NDM, KPC, OXA, IMP, VIM, IMI and GES classes were identified, the most
250 frequent being NDM-5 (n=893), KPC-2 (n=264), NDM-1 (n=219) and OXA-48 (n=189).

251

252 We generated model estimated mean proportions with 95% CIs for models of MDR (Fig. 2),
253 XDR (Fig. S3) and *bla*_{CTX-M} (Fig. 3) presence, according to host, phylogroup and subregion
254 for both the full dataset, and a dataset subsampled down to the smallest group size without
255 replacement to adjust for large sample size differences between groups. The same process
256 was repeated for carbapenemase gene presence (Fig. 4), but subregion was not used as a
257 predictor due the rarity of carbapenemase-positive genomes across many regions. Trends
258 were broadly consistent across the two datasets and estimates for the full dataset are quoted in
259 the subsequent text.

260

261 *1. MDR and XDR*

262

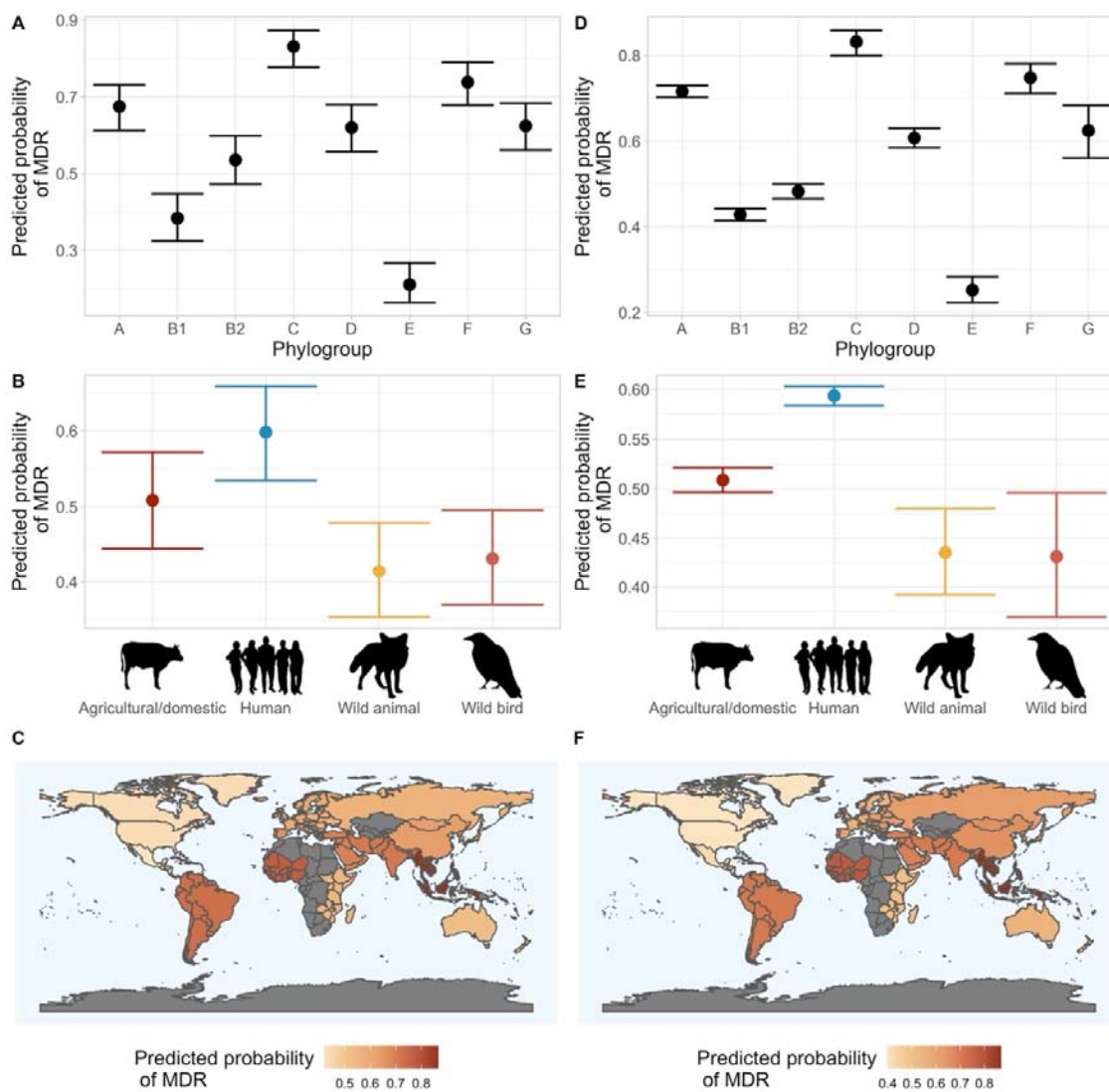
263 MDR genomes were common, with a given genome from the majority of phylogroups,
264 subregions and hosts having > 0.4 probability of being MDR (Fig. 2). Within phylogroups,
265 multidrug resistance was most common in groups C (0.83, 95% CI=0.80-0.86), F (0.75, 95%
266 CI=0.71-0.78) and A (0.72, 95% CI=0.70-0.73), whilst being comparatively uncommon in
267 group E isolates (0.25, 95% CI=0.22-0.28). Isolates from humans (0.60, 95% CI=0.58-0.60)
268 and agricultural/domestic hosts (0.51, 95% CI=0.50-0.52) had the highest probabilities of

269 MDR, though estimates from wild animals (0.44, 95% CI=0.39-0.48) and birds (0.43, 95%
270 CI= 0.37-0.50) were reasonably high. Finally, South-Eastern Asia (0.86, 95% CI = 0.83-0.90)
271 and Western Africa (0.78, 95% CI=0.15-0.72) had the highest estimated proportions of MDR
272 genomes, while the lowest were in Northern America (0.40, 95% CI=0.38-0.41) and Central
273 America (0.42, 95% CI=0.33-0.52).

274

275 XDR genomes were less common, and not detected at high enough frequencies to estimate
276 their probabilities in the minimum subsampled dataset for phylogroup B2 and wild animal
277 hosts (Fig. S3). The phylogroups with the highest probabilities of XDR were A (0.092, 95%
278 CI=0.083-0.10) and F (0.07, 95% CI=0.052-0.093), whilst remaining very low for B2 (0.003,
279 95% CI=0.0018-0.0058) and E (0.018, 95% CI=0.012-0.03). The highest XDR probabilities
280 were estimated in wild birds (0.13, 95% CI=0.091-0.18), with probabilities below 0.06 for
281 other host categories. Finally, XDR genomes were by far the most likely to occur in Eastern
282 Asia (0.16, 95% CI=0.14-0.17), and below 0.08 for all other regions.

283



284

285 Figure 2 –Estimated mean probabilities and 95% CIs of MDR from binomial GLMs for datasets
 286 subsampled without replacement to the minimum group size for A) phylogroup, B) host category and
 287 C) geographic subregion. The same model types for the full non-subsampled dataset for D)
 288 phylogroup, E) host category and F) geographic subregion. CIs for subregion models can be found in
 289 Figure S6.

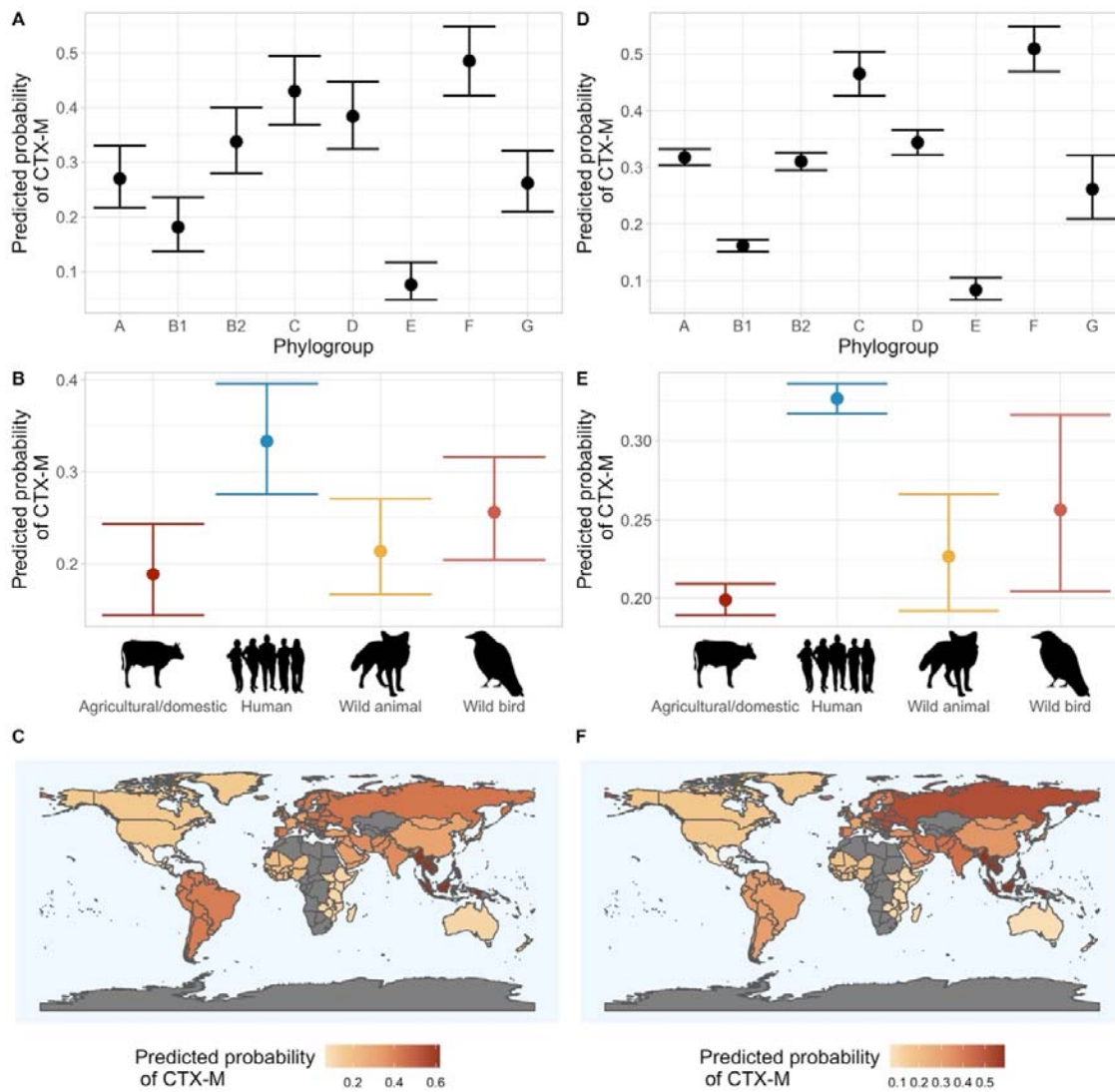
290

291 2. *bla_{CTX-M}* and *carbapenemase* genes

292

293 The presence of the *bla_{CTX-M}* class of β -lactamases was highly variable between groups (Fig.
 294 3). Estimated probabilities ranged from 0.08 (95% CI = 0.066-0.11) in phylogroup E to 0.51
 295 (95% CI=0.47-0.55) in phylogroup F, whilst *bla_{CTX-M}* genes were more associated with
 296 human hosts (0.33, 95% CI=0.32-0.34) than other host categories. The probabilities of

297 *bla*_{CTX-M} genes being present also varied widely between subregions, with the highest in
 298 South-Eastern Asia (0.59, 95% CI=0.53-0.64) and Eastern Europe (0.52, 95% CI=0.47-0.57),
 299 and all other subregions below 0.45.



300
 301 Figure 3 – Estimated mean probabilities and 95% CIs of *bla*_{CTX-M} presence from binomial GLMs for
 302 datasets subsampled without replacement to the minimum group size for A) phylogroup, B) host
 303 category and C) geographic subregion. The same model types for the full non-subsampled dataset for
 304 D) phylogroup, E) host category and F) geographic subregion. CIs for subregion models can be found
 305 in Figure S12.

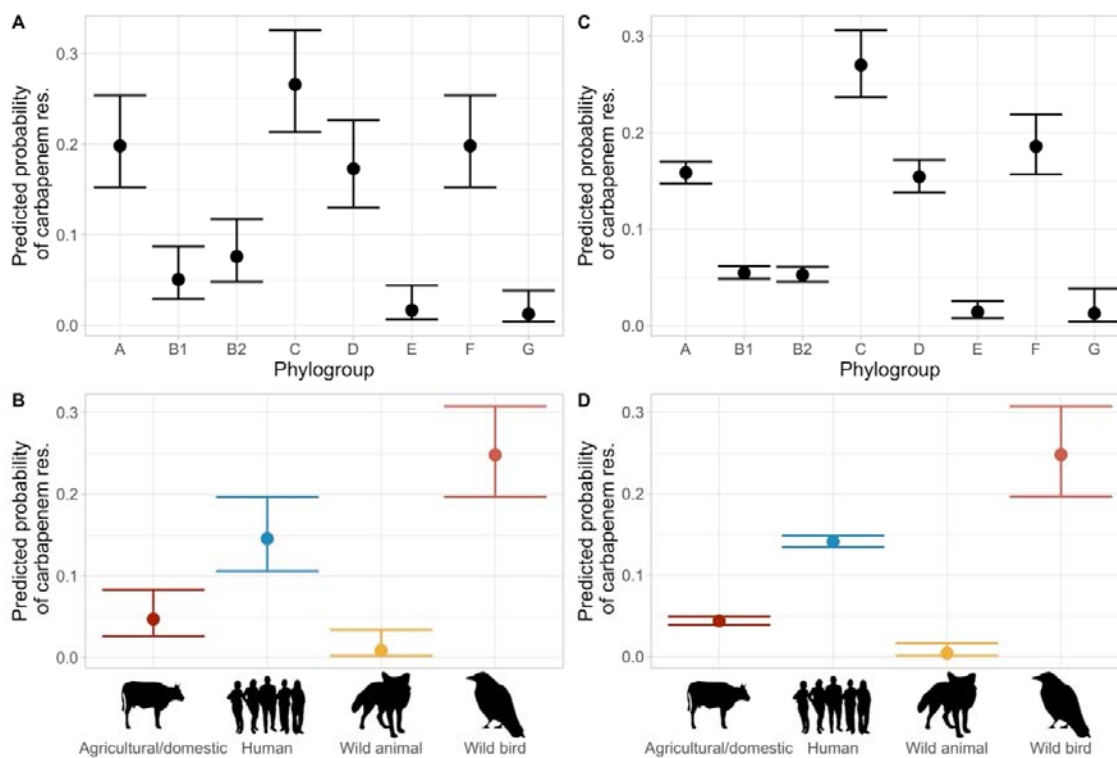
306
 307 Carbapenemase genes were far less common than *bla*_{CTX-M} genes (Fig. 4). There was a clear
 308 split between phylogroups C, F, A and D, which were more likely to possess carbapenemases
 309 (probabilities >0.15), and the remaining groups which had probabilities of 0.05 or below. The

310 highest proportions of carbapenemases were estimated in wild bird (0.25, 95% CI=0.20-0.31)
 311 and human (0.14, 95% CI= 0.13-0.15) hosts.

312

313 We further investigated sampling bias for carbapenemase-positive isolates in wild birds due
 314 to the high proportion of genomes possessing these genes. Isolates from wild birds possessing
 315 carbapenemases were only represented in 3 subregions: Australia and New Zealand (n=25),
 316 Eastern Asia (n=32) and Western Asia (n=1). Sequencing bias was evident; for example, 68
 317 out of 234 wild bird isolates were from one study of silver gulls in Australia (which made up
 318 68 out of the 93 isolates from this region). This study (PRJNA630096) specifically sought
 319 and sequenced isolates conferring resistance to critically important antimicrobials. Further
 320 carbapenemase-producing isolates came from studies in China (PRJNA669620 and
 321 PRJNA349231) that specifically studied carbapenem-resistant *E. coli*.

322



323

324 Figure 4 – Estimated mean probabilities and 95% CIs of carbapenemase presence from binomial
 325 GLMs for datasets subsampled without replacement to the minimum group size for A) phylogroup
 326 and B) host category. The same model types for the full non-subsampled dataset for C) phylogroup
 327 and D) host category.

328

329 *Subsampling and resampling approaches do not alter major trends*

330

331 Due to the large discrepancies in sample size between hosts, regions and phylogroups in this
332 dataset, we used additional sampling with replacement to examine the sensitivity of the
333 observed trends to changing sample sizes of 500, 1000, 2000 and 4000 (Fig. S4-S12). This
334 showed that, although fluctuations in estimated proportions occurred (as well as expected
335 reductions in CIs with increased sample size), broader trends were reproducible across wide
336 variations in subsample size. In some smaller subsamples, insufficient genomes possessed the
337 metric of AMR carriage (e.g. a carbapenemase gene) to make estimates in some
338 subgroupings.

339

340 **Discussion**

341

342 Though we are in the era of large-scale genomics studies, with thousands of genomes
343 available for species such as *E. coli*, it remains challenging to collate this information to
344 answer questions about AMR prevalence. In this study, we characterised the RefSeq *E. coli*
345 dataset for which metadata was available, investigating sampling bias and AMR carriage
346 according to various metrics for different hosts, phylogroups and geographic subregions.

347

348 *How evenly sampled are publicly-available E. coli genomes?*

349

350 In first characterising this dataset, large sample size discrepancies in the locations and hosts
351 *E. coli* genomes were sourced from were apparent. This was not unexpected. Previous large-
352 scale work on *E. coli* genomes isolated from human hosts has emphasised the sampling
353 biases towards clinically-important isolates and lineages, as well as high-income countries,
354 even in the most well-sampled host species [13]. Other work has used literature searches to
355 quantify how poorly sampled wild animals are for *E. coli*, with small within-group sample
356 sizes in the few studies available [28]. Our data adds to this body of knowledge emphasising
357 how little is known about ARG burden across the diversity of this species, despite the vast
358 amount of sequencing data available.

359

360 We investigated the effects of subsampling all groups to the minimum sample size, as well as
361 resampling up to larger sample sizes. This approach cannot replace unbiased sampling,
362 particularly for those groups that are less well represented. However, it is interesting to note
363 that estimates of AMR carriage for larger groups, such as phylogroup A or human hosts, did

364 not alter substantially regardless of the subsampling regime. Additionally, the spread of
365 counts of ARGs was wide in these groups, indicating that they were not necessarily
366 dominated by one or few lineages.

367

368 Previous work has also investigated the potential effect of sampling bias on AMR estimates
369 using *E. coli* genomes from GenBank [29]. These authors found that 1) the average distance
370 of a given newly sequenced *E. coli* genome is equivalent to that between known close
371 relatives, such as O157:H7 genomes, and 2) the results of a population genetics analysis of
372 AMR did not change substantially when all human or unknown source genomes were
373 removed, or all genomes for which the GenBank record referenced antibiotic resistance.
374 Therefore, it is possible that, for more well-sampled subgroupings of *E. coli*, our current
375 sequence collection already spans a high diversity of this species.

376

377 *Phylogroup distributions within hosts and subregions*

378

379 Though the most common phylogroups were A, B1 and B2, their proportions varied
380 substantially between hosts. Phylogroup B1 was the most frequent in agricultural and wild
381 animals. This is broadly consistent with many smaller-scale studies for which B1 was the
382 most frequent in agricultural settings. For example, B1 made up 71% of agriculture-
383 associated isolates in a study in the Philippines [30], 50% of those in a study of sheep farming
384 in China [31], and 35% of isolates from cattle and their attendants in Tanzania [32]. The wide
385 variation in actual proportions observed in these different studies may be due to smaller
386 sample sizes (n=17-100 in those referenced), as well as location and species effects. Human
387 isolates were proportionally more likely to be B2 than those from other hosts. *E. coli* isolates
388 from group B2 tend to possess more virulence traits than those from A and B1, though this
389 virulence is proposed to be a by-product of commensalism [33]. Geographic differences in
390 predominant commensal *E. coli* phylogroups have broadly represented a shift from A to B2
391 in industrialised nations over the past decades, with these being the most common
392 phylogroups in humans [34]. We saw this trend to some extent – European regions and
393 Australia and New Zealand had high proportions of B2. Contrastingly, we found that other
394 regions such as North America and Eastern Asia had B1 as the predominant phylogroup.

395

396 *Estimating AMR carriage across phylogroups, hosts and subregions*

397

398 We used various metrics of AMR carriage to get an overview of ARG distribution in
399 publicly-available *E. coli* genomes. When comparing phylogroups, A, C and F were
400 repeatedly associated with more ARGs as well as clinically-important β -lactamases.
401 Phylogroup A is a major group that appears to be more generalist. It is spread across
402 vertebrate hosts [8], and contains mostly human commensal strains [35] as well as laboratory
403 strains [36]. Although the presence of engineered strains may partly explain the high AMR
404 metrics in this group, our results also imply that ARGs may be frequent in commensal
405 lineages. Phylogroups commonly present in the human gut could be a reservoir for resistance
406 genes and plasmids [37], potentially allowing their horizontal spread into other, more
407 problematic lineages or species.

408

409 As they are minor phylogroups, comparatively less is known about groups C and F. In terms
410 of pathogenic potential, the sequence type complex 88 (STc88) lineage of group C is one of
411 the main avian pathogenic *E. coli* (APEC) strains [38]. This phylogroup also includes ST410,
412 a recently-emerged MDR lineage with fluoroquinolone resistance, ESBL (*bla*_{CTX-M-15}) and
413 carbapenemase (*bla*_{OXA-181} and *bla*_{NDM-5}) genes [38]. Meanwhile, phylogroup F has been
414 associated with fewer virulence traits than closely related groups [39], though it has also been
415 linked to extraintestinal pathogenic (ExPEC) infections [38] and contains the MDR STc648
416 lineage [40]. Therefore, these phylogroups could represent highly virulent and/or drug-
417 resistant lineages and may warrant further investigation.

418

419 Unsurprisingly, human hosts were associated with more ARGs than other host categories.
420 However, high levels were also seen in agricultural and domestic animal species for some
421 measures, such as MDR. In addition, *bla*_{CTX-M} and carbapenemase genes were detected in all
422 the animal host categories to some degree. Antibiotics are still extensively used in animal
423 agriculture, selecting for ARGs both directly in the gut and through excretion of antibiotics
424 into soil and the local environment [41]. Despite this there is no standardised global
425 surveillance system used in animal agriculture that is equivalent to those used in humans
426 [41]. There is evidence for transmission of *bla*_{CTX-M} variants in *E. coli* from humans to farmed
427 animals and the environment, which also occurs in the reverse direction via food products
428 [42]. Only limited studies document the occurrence of carbapenemases in *E. coli* strains from
429 wildlife, agricultural animals and soils, though they have been reported in these contexts [43,
430 44]. This highlights the need for further prevalence studies to determine the extent to which
431 these genes are spreading in different environments.

432

433 The proportionally high numbers of carbapenemase-possessing isolates (and potentially XDR
434 isolates) from wild birds in this work are likely due to studies that enriched for ARGs.

435 However, this emphasizes our lack of knowledge on the true extent of AMR spread in the
436 natural world. High levels of problematic ARGs in wild birds, particularly generalist species,
437 could occur due to their relatively high mobility compared to other taxa, and their utilisation
438 of diverse foraging habitats, both of which increase potential exposure to anthropogenic
439 sources of resistant bacteria [45, 46].

440

441 Finally, estimates of ARG prevalence varied geographically depending on the metric used.
442 The highest AMR carriage was generally estimated for Eastern and South-Eastern Asia
443 across measures. However, other regions were particularly high for individual measures. For
444 example, Eastern Europe had high levels of *bla*_{CTX-M} genes and MDR genomes were more
445 common in Western Africa. These trends could either reflect genuine differences in AMR
446 levels and individual gene prevalence, or regional differences in sampling strategies and
447 genomic surveillance. There are large gaps in AMR surveillance for low and middle income
448 (LMIC) countries. A recent study utilised the relationship between socioeconomic
449 characteristics and AMR prevalence to model AMR levels for underrepresented countries,
450 estimating high levels of third-generation cephalosporin resistance in *E. coli* from Western
451 Asia [47]. Genomic surveillance has been leveraged during the COVID-19 pandemic and for
452 other infectious diseases, but next-generation sequencing capacity is low in many regions
453 [48]. However, the development of more many sequence analysis tools to detect ARGs [49],
454 as well as developing capacity to incorporate whole genome sequencing into AMR
455 surveillance in LMICs [50], may lead to this becoming a powerful tool in the future.

456

457 *Public genomic databases – limitations, benefits and future work*

458

459 We show the void in sequencing data from *E. coli* isolates outside of human and agricultural
460 settings in a very large dataset, as well as stark regional disparities, with scarce publicly-
461 available data representing the majority of global regions. Whilst this study characterises this
462 currently available dataset, representative estimates of ARG burden in *E. coli* sequences
463 cannot be made without systematic, representative sampling. Notably, it is concerning that
464 ARGs that inactivate our last-resort antibiotics are being detected and potentially spread in
465 any environmental context, and that we do not have sufficient genomic data to investigate

466 this. Despite their limitations, publicly-available datasets are important, allowing access to
467 sequencing data for scientists and the public across the globe, as well as improving
468 reproducibility. Future systematic sampling and genomic surveillance in previously under-
469 investigated settings will give a more accurate picture of the scale of the AMR problem.

470

471 **Data Summary**

472

473 All code used to generate and analyse data are available at
474 <https://github.com/elliekpurse/AMR-Ecoli>.

475

476 **Author contributions**

477

478 EP, SVH and ERW conceptualised the study. WHG refined the directions of the analysis and
479 methodology. EP performed all data and statistical analysis and wrote the first manuscript
480 draft. All authors contributed to interpretation of the results and editing subsequent versions
481 of the manuscript.

482

483 **Conflicts of interest**

484

485 We declare no conflicts of interest.

486

487 **Acknowledgements and funding**

488

489 EP was supported by a PhD studentship equally funded by a grant from the European
490 Research Council under the European Union's Horizon 2020 research and innovation
491 programme (ERC-STG-2016-714478 to ERW) and the College of Life and Environmental
492 Sciences, University of Exeter, UK. EP thanks Liam Langley for statistical modelling advice,
493 and Elze Hesse and Kate Baker for feedback on the analyses. TD is supported by European
494 Research Council (ERC-2017-ADG-788405 to ERW). SVH acknowledges BBSRC for
495 funding (BB/R010781/1 and BB/S017674/1).

496

497 **References**

498

- 499 1. **Luo C, Walk ST, Gordon DM, Feldgarden M, Tiedje JM, et al.** Genome sequencing of
500 environmental *Escherichia coli* expands understanding of the ecology and speciation of
501 the model bacterial species. *Proc Natl Acad Sci* 2011;108:7200–7205.
- 502 2. **Murray CJ, Ikuta KS, Sharara F, Swetschinski L, Aguilar GR, et al.** Global burden of
503 bacterial antimicrobial resistance in 2019: a systematic analysis. *The Lancet*
504 2022;399:629–655.
- 505 3. **World Health Organization.** *Prioritization of pathogens to guide discovery, research and*
506 *development of new antibiotics for drug-resistant bacterial infections, including*
507 *tuberculosis.* 9789240026438 (electronic version); Geneva: World Health Organization;
508 2017.
- 509 4. **Hall RJ, Whelan FJ, Cummins EA, Connor C, McNally A, et al.** Gene-gene
510 relationships in an *Escherichia coli* accessory genome are linked to function and mobility.
511 *Microb Genomics* 2021;7:000650.
- 512 5. **Von Wintersdorff CJH, Penders J, Van Niekerk JM, Mills ND, Majumder S, et al.**
513 Dissemination of antimicrobial resistance in microbial ecosystems through horizontal
514 gene transfer. *Front Microbiol* 2016;7:1–10.
- 515 6. **Zhang A-N, Gaston JM, Dai CL, Zhao S, Poyet M, et al.** An omics-based framework
516 for assessing the health risk of antimicrobial resistance genes. *Nat Commun*
517 2021;12:4765.
- 518 7. **Woolhouse M, Ward M, van Bunnik B, Farrar J.** Antimicrobial resistance in humans,
519 livestock and the wider environment. *Philos Trans R Soc B Biol Sci* 2015;370:20140083–
520 20140083.
- 521 8. **Touchon M, Perrin A, Sousa JAM de, Vangchhia B, Burn S, et al.** Phylogenetic
522 background and habitat drive the genetic diversification of *Escherichia coli*. *PLOS Genet*
523 2020;16:e1008866.
- 524 9. **Economou V, Gousia P.** Agriculture and food animals as a source of antimicrobial-
525 resistant bacteria. *Infect Drug Resist* 2015;8:49–61.
- 526 10. **Arnold KE, Williams NJ, Bennett M.** ‘Disperse abroad in the land’: the role of wildlife
527 in the dissemination of antimicrobial resistance. *Biol Lett* 2016;12:20160137.
- 528 11. **Robinson TP, Bu DP, Carrique-Mas J, Fèvre EM, Gilbert M, et al.** Antibiotic
529 resistance is the quintessential One Health issue. *Trans R Soc Trop Med Hyg*
530 2016;110:377–380.
- 531 12. **Petitjean M, Condamine B, Burdet C, Denamur E, Ruppé E.** Phylum barrier and
532 *Escherichia coli* intra-species phylogeny drive the acquisition of antibiotic-resistance
533 genes. *Microb Genomics* 2021;7:000489.
- 534 13. **Horesh G, Blackwell GA, Tonkin-Hill G, Corander J, Heinz E, et al.** A
535 comprehensive and high-quality collection of *Escherichia coli* genomes and their genes.
536 *Microb Genomics* 2021;7:000499.

- 537 14. **Feldgarden M, Brover V, Gonzalez-Escalona N, Frye JG, Haendiges J, et al.**
538 AMRFinderPlus and the Reference Gene Catalog facilitate examination of the genomic
539 links among antimicrobial resistance, stress response, and virulence. *Sci Rep*
540 2021;11:12728.
- 541 15. **Clermont O, Christenson JK, Denamur E, Gordon DM.** The Clermont *Escherichia*
542 *coli* phylo-typing method revisited: improvement of specificity and detection of new
543 phylo-groups. *Environ Microbiol Rep* 2013;5:58–65.
- 544 16. **Waters NR, Abram F, Brennan F, Holmes A, Pritchard L.** Easy phylotyping of
545 *Escherichia coli* via the EzClermont web app and command-line tool. *Access Microbiol.*
- 546 17. **Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, et al.** *Sustainable*
547 *data analysis with Snakemake.* 10:33; F1000Research; 18 January 2021.
- 548 18. **Magiorakos A-P, Srinivasan A, Carey RB, Carmeli Y, Falagas ME, et al.** Multidrug-
549 resistant, extensively drug-resistant and pandrug-resistant bacteria: an international expert
550 proposal for interim standard definitions for acquired resistance. *Clin Microbiol Infect*
551 2012;18:268–281.
- 552 19. **Skarżyńska M, Zając M, Bomba A, Bocian Ł, Kozdruń W, et al.** Antimicrobial
553 Resistance Glides in the Sky—Free-Living Birds as a Reservoir of Resistant *Escherichia*
554 *coli* With Zoonotic Potential. *Front Microbiol*;12.
555 <https://www.frontiersin.org/articles/10.3389/fmicb.2021.656223> (2021, accessed 4 July
556 2022).
- 557 20. **Wickham H, Averick M, Bryan J, Chang W, McGowan L, et al.** Welcome to the
558 Tidyverse. *J Open Source Softw* 2019;4:1686.
- 559 21. **Mills BR.** MetBrewer: Color Palettes Inspired by Works at the Metropolitan Museum of
560 Art.
- 561 22. **Pebesma E.** Simple Features for R: Standardized Support for Spatial Vector Data. *R J*
562 2018;10:439.
- 563 23. **Brooks M E, Kristensen K, Benthem K J, van, Magnusson A, Berg C W, et al.**
564 glmmTMB Balances Speed and Flexibility Among Packages for Zero-inflated
565 Generalized Linear Mixed Modeling. *R J* 2017;9:378.
- 566 24. **Hartig F.** DHARMA: Residual Diagnostics for Hierarchical (Multi-Level / Mixed)
567 Regression Models. <http://florianhartig.github.io/DHARMA/> (2020).
- 568 25. **Lüdtke D.**ggeffects: Tidy Data Frames of Marginal Effects from Regression Models. *J*
569 *Open Source Softw* 2018;3:772.
- 570 26. **Sullivan GM, Feinn R.** Using Effect Size—or Why the P Value Is Not Enough. *J Grad*
571 *Med Educ* 2012;4:279–282.
- 572 27. **Berner D, Amrhein V.** Why and how we should join the shift from significance testing
573 to estimation. *J Evol Biol*;n/a.

- 574 28. **Lagerstrom KM, Hadly EA.** The under-investigated wild side of *Escherichia coli*:
575 genetic diversity, pathogenicity and antimicrobial resistance in wild animals. *Proc R Soc*
576 *B Biol Sci* 2021;288:20210399.
- 577 29. **Goldstone RJ, Smith DGE.** A population genomics approach to exploiting the accessory
578 'resistome' of *Escherichia coli*. *Microb Genomics* 2017;3:e000108.
- 579 30. **Zara ES, Vital PG.** Phylogroup typing and carbapenem resistance of *Escherichia coli*
580 from agricultural samples in Metro Manila, Philippines. *J Environ Sci Health Part B*
581 2022;57:644–656.
- 582 31. **Zhao X, Lv Y, Adam FEA, Xie Q, Wang B, et al.** Comparison of Antimicrobial
583 Resistance, Virulence Genes, Phylogroups, and Biofilm Formation of *Escherichia coli*
584 Isolated From Intensive Farming and Free-Range Sheep. *Front Microbiol*;12.
585 <https://www.frontiersin.org/articles/10.3389/fmicb.2021.699927> (2021, accessed 19 July
586 2022).
- 587 32. **Madoshi BP, Kudirkiene E, Mtambo MMA, Muhairwa AP, Lupindu AM, et al.**
588 Characterisation of Commensal *Escherichia coli* Isolated from Apparently Healthy Cattle
589 and Their Attendants in Tanzania. *PLOS ONE* 2016;11:e0168160.
- 590 33. **Le Gall T, Clermont O, Gouriou S, Picard B, Nassif X, et al.** Extraintestinal Virulence
591 Is a Coincidental By-Product of Commensalism in B2 Phylogenetic Group *Escherichia*
592 *coli* Strains. *Mol Biol Evol* 2007;24:2373–2384.
- 593 34. **Tenaillon O, Skurnik D, Picard B, Denamur E.** The population genetics of commensal
594 *Escherichia coli*. *Nat Rev Microbiol* 2010;8:207–217.
- 595 35. **Smati M, Clermont O, Le Gal F, Schichmanoff O, Jauréguy F, et al.** Real-Time PCR
596 for Quantitative Analysis of Human Commensal *Escherichia coli* Populations Reveals a
597 High Frequency of Subdominant Phylogroups. *Appl Environ Microbiol* 2013;79:5005–
598 5012.
- 599 36. **Abram K, Udaondo Z, Bleker C, Wanchai V, Wassenaar TM, et al.** Mash-based
600 analyses of *Escherichia coli* genomes reveal 14 distinct phylogroups. *Commun Biol*
601 2021;4:1–12.
- 602 37. **Hu Y, Yang X, Qin J, Lu N, Cheng G, et al.** Metagenome-wide analysis of antibiotic
603 resistance genes in a large cohort of human gut microbiota. *Nat Commun* 2013;4:2151.
- 604 38. **Denamur E, Clermont O, Bonacorsi S, Gordon D.** The population genetics of
605 pathogenic *Escherichia coli*. *Nat Rev Microbiol*. Epub ahead of print 21 August 2020.
606 DOI: 10.1038/s41579-020-0416-x.
- 607 39. **Vangchhia B, Abraham S, Bell JM, Collignon P, Gibson JS, et al.** Phylogenetic
608 diversity, antimicrobial susceptibility and virulence characteristics of phylogroup F
609 *Escherichia coli* in Australia. *Microbiology* 2016;162:1904–1912.
- 610 40. **Johnson JR, Johnston BD, Gordon DM.** Rapid and Specific Detection of the
611 *Escherichia coli* Sequence Type 648 Complex within Phylogroup F. *J Clin Microbiol*
612 2017;55:1116–1121.

- 613 41. **Jadeja NB, Worrich A.** From gut to mud: dissemination of antimicrobial resistance
614 between animal and agricultural niches. *Environ Microbiol* 2022;24:3290–3306.
- 615 42. **Bevan ER, Jones AM, Hawkey PM.** Global epidemiology of CTX-M β -lactamases:
616 temporal and geographical shifts in genotype. *J Antimicrob Chemother* 2017;72:2145–
617 2155.
- 618 43. **Furlan JPR, Sellera FP, Stehling EG.** Strengthening genomic surveillance of
619 carbapenemases in soils: a call for global attention. *Lancet Microbe*;0. Epub ahead of
620 print 23 March 2023. DOI: 10.1016/S2666-5247(23)00093-9.
- 621 44. **Köck R, Daniels-Haardt I, Becker K, Mellmann A, Friedrich AW, et al.**
622 Carbapenem-resistant Enterobacteriaceae in wildlife, food-producing, and companion
623 animals: a systematic review. *Clin Microbiol Infect* 2018;24:1241–1250.
- 624 45. **Navarro J, Grémillet D, Afán I, Miranda F, Bouten W, et al.** Pathogen transmission
625 risk by opportunistic gulls moving across human landscapes. *Sci Rep* 2019;9:10659.
- 626 46. **Migura-Garcia L, Ramos R, Cerdà-Cuéllar M.** Antimicrobial Resistance of
627 *Salmonella* Serovars and *Campylobacter* spp. Isolated from an Opportunistic Gull
628 Species, Yellow-legged Gull (*Larus michahellis*). *J Wildl Dis* 2017;53:148–152.
- 629 47. **Oldenkamp R, Schultsz C, Mancini E, Cappuccio A.** Filling the gaps in the global
630 prevalence map of clinical antimicrobial resistance. *Proc Natl Acad Sci*
631 2021;118:e2013515118.
- 632 48. **Inzaule SC, Tessema SK, Kebede Y, Ogwel Ouma AE, Nkengasong JN.** Genomic-
633 informed pathogen surveillance in Africa: opportunities and challenges. *Lancet Infect Dis*
634 2021;21:e281–e289.
- 635 49. **Hendriksen RS, Bortolaia V, Tate H, Tyson GH, Aarestrup FM, et al.** Using
636 Genomics to Track Global Antimicrobial Resistance. *Front Public Health*;7.
637 <https://www.frontiersin.org/articles/10.3389/fpubh.2019.00242> (2019, accessed 22 May
638 2023).
- 639 50. **NIHR Global Health Research Unit on Genomic Surveillance of AMR.** Whole-
640 genome sequencing as part of national and international surveillance programmes for
641 antimicrobial resistance: a roadmap. *BMJ Glob Health* 2020;5:e002244.

642