

1 **Appendix for:**  
2 **Multidrug Resistance Dynamics in *Salmonella* in Food-Animals in the United**  
3 **States: An Analysis of Genomes from Public Databases**

4 João Pires<sup>1\*</sup>, Jana S. Huisman<sup>2,3</sup>, Sebastian Bonhoeffer<sup>2</sup>, and Thomas P. Van Boeckel<sup>1,4</sup>

5

6 **Affiliations:**

7

8 1. Institute for Environmental Decisions, ETH Zurich, Zurich, Switzerland

9 2. Institute of Integrative Biology, ETH Zurich, Zurich, Switzerland

10 3. Swiss Institute of Bioinformatics, Lausanne, Switzerland

11 4. Center for Disease Dynamics, Economics & Policy, New Delhi, India

12

13

14

15 This appendix contains:

16 I. Supplementary Materials and Methods

17 II. Appendix Figures

18 III. Legends for Appendix Tables

19 IV. Legend for Appendix Dataset

20

21

22

23

24

25

26

27

28 I. Supplementary Materials and Methods

29

30 1. Data Retrieval and Harmonization

31 1.1. Genomic Metadata Retrieval

32

33 We searched for all available *Salmonella enterica* assemblies recovered from food-animals  
34 released until the end of 2018 in three public genomic data repositories: the National Center  
35 for Biotechnology Information (NCBI) Nucleotide database(1), EnteroBase(2) and  
36 Pathosystems Resource Integration Center (PATRIC)(3). For NCBI Nucleotide, we queried  
37 Entrez Programming Utilities using the taxonomic identification of *S. enterica*  
38 (“taxid28901”)(4). Resulting accession numbers were retrieved and used to retrieve  
39 associated metadata. For both EnteroBase and PATRIC, the entire metadata tables were  
40 downloaded.

41

42 1.2. Metadata Standardization

43

44 Metadata tables were imported into R (version 3.6.0)(5). Manipulation of metadata was  
45 performed with the tidyverse(6) (version 1.3.0), data.table(7) (version 1.12.8) and plyr(8)  
46 (version 1.8.6) packages.

47 All entries that did not report a country of origin or geographic coordinates were removed.  
48 Thereafter, we inspected isolation sources and to identify food-animal key words that allowed  
49 to reduce the datasets, but would maximize the number of hits. The resulting filtered datasets  
50 were then manually curated to exclude entries that did not meet the criteria of food-animal.  
51 We considered four levels of data aggregation for host attribution:

- 52 • **Source Niche:** highest level of aggregation and indicates whether samples were  
53 recovered from food-products (Food) or from the animals themselves (Poultry or  
54 Livestock);
- 55 • **Generic Host:** aggregation within animal-production group such as poultry, swine,  
56 bovine, ovine and caprine. The categories dairy, meat and environment were also  
57 introduced when no specific animal was given. The category environment denotes  
58 food-animal-related samples not collected directly from the animal or their food-  
59 products such as drag swabs, poultry litter, eggshells, animal bedding and barns;
- 60 • **Source Type:** indicates the specific animal from which the samples were collected;
- 61 • **Source Details:** contains the original sample description as input by the submitter.

62

63 Geographic coordinates were retrieved from metadata tables when available. Coordinates  
64 expressed in cardinal directions were converted to decimal degree. For entries without  
65 coordinates, an address was constructed based on the available information of the isolation  
66 location (country, province, state, region, city, zip code, etc.). We queried addresses for their  
67 decimal degree coordinates with the geocode function of the ggmap package(9) (version  
68 3.0.0). Assemblies returning no coordinates were inspected manually and queried in Google  
69 Earth Pro(10). A column with country's three-letter code based on the ISO 3166-1 guidelines  
70 was also assigned.

71 Isolation dates were harmonized according to the ISO 8601 format (year-month-day) using  
72 lubridate(11) and anytime(12) packages. A dedicated column for year of isolation was also  
73 created. Finally, the NCBI BioSample and BioProject (when available) were also kept.

74

### 75 1.3. Creation of a Consensus Dataset and Assembly Download 76

77 We used the BioSample identifier to compare entries across databases and created a  
78 consensus dataset by removing duplicate entries. We primarily kept entries from EnteroBase  
79 given the dedicated pipeline this database has towards short-read sequences assembly,  
80 quality control, and molecular typing). Then we retrieved data from PATRIC and finally from  
81 NCBI RefSeq(13). EnteroBase derived assemblies were kindly provided by the curators,  
82 PATRIC assemblies were downloaded through the PATRIC Command Line Interface(14), NCBI  
83 assemblies were downloaded from the RefSeq database(13). Each entry has the original  
84 identifier from the database and a column indicating from which database it retrieved from.

85

### 86 2. Curation of Predicted Phenotypes 87

88 We extracted the predicted phenotypes from ResFinder database  
89 ([https://bitbucket.org/genomicepidemiology/resfinder\\_db/src/master/](https://bitbucket.org/genomicepidemiology/resfinder_db/src/master/), accessed 27<sup>th</sup> May  
90 2020). We retrieved the predicted antibiotic family (Antibiotic Class) and specific antibiotics  
91 (Phenotype) to which they confer resistance. All antimicrobial resistance genes (ARGs) found  
92 in our dataset can be found in Supplementary Table 2. Phenotypes of genes with unassigned  
93 predicted phenotypes were inputted based on the closest match sequence match. In brief,  
94 we retrieved the sequence of such ARGs based on the available NCBI accession number and  
95 used Basic Local Alignment Search Tool (BLAST)(15) against the CARD database. Predicted  
96 phenotypes were assigned based on the gene with the best alignment score, but with a  
97 minimum of 97% identity. When no matches were found in CARD, we used NCBI's BLAST(16)  
98 instead. For the matches with the highest identity and coverage, we inspected the referred

99 manuscripts where such ARG and their respective resistance phenotypes were described.  
100 Finally, for  $\beta$ -lactamase genes, we added cephalothin manually to the Phenotype column. This  
101 is because early generation cephalosporins were not included in the ResFinder phenotype list,  
102 although TEM-types(17), AmpC(18), OXA-Types(19) hydrolyze these  $\beta$ -lactams. We removed  
103 *aac(6')-Iaa* from the dataset as this gene has been described as intrinsic to *S. enterica* and  
104 does not cause phenotypic resistance(20, 21).

105 In the case of point mutations, we only kept those that have known resistance phenotype in  
106 the PointFinder database (22), which can be extracted directly from the staramr output(23).

107

### 108 3. Multidrug Resistance Score Calculation

109

110 We devised a metric to summarize multidrug resistance based on the number of different  
111 classes of antimicrobials an isolate is predicted to have resistance based on its content in  
112 ARGs (acquired ARGs or point mutations). We call this metric the Multidrug Resistance Score  
113 (MDR Score). We based this metric on microbiological resistance. In brief, microbiological  
114 resistance is identified when the minimal inhibitory concentration (MIC) is above the  
115 epidemiological cut-off (ECOFF) value(24) – the highest MIC for organisms devoid of  
116 phenotypically detectable acquired resistance mechanisms(25). We assume that all identified  
117 ARGs are functional and thus resulting MICs would be above the ECOFF. For the majority of  
118 ARGs, the phenotype would not be affected. However, it could affect the predicted  
119 phenotype  $\beta$ -lactamase genes since for some variants the amino acid changes result in  
120 different resistance phenotypes(26). Although only 0.44% of the  $\beta$ -lactamase genes had a  
121 coverage and identity below 100% in our study.

122 To calculate the MDR Score, we used a list of antimicrobials of clinical importance  
123 *Enterobacteriaceae* in relation to acquired resistance(27). We used cephalothin as a surrogate  
124 for cefazolin since they are both early generation cephalosporins.

125 The MDR score was computed as follow:

126 • For each genome, the unique predicted resistance phenotypes were identified, and  
127 antibiotics were grouped into the different molecular classes:

- 128 ○ Aminoglycosides
- 129 ○ Penicillins
- 130 ○ Early Generation Cephalosporins
- 131 ○ Cephamycins
- 132 ○ 3<sup>rd</sup> Generation Cephalosporins
- 133 ○ 4<sup>th</sup> Generation Cephalosporins
- 134 ○ Monobactams
- 135 ○ Carbapenems
- 136 ○ Penicillins in combination with  $\beta$ -lactamase inhibitors
- 137 ○ Quinolones
- 138 ○ Trimethoprim
- 139 ○ Sulphonamides
- 140 ○ Phenicol
- 141 ○ Tetracyclines
- 142 ○ Polymyxins
- 143 ○ Fosfomicin

144

145 • The MDR Score will increase by one when an antibiotic is assigned to one of the  
146 described molecular classes. If more ARGs confer resistance to the same molecular  
147 class, the MDR score still only increases by one.

148 • All genomes for which no ARGs are identified are assigned a MDR score of zero.

149

#### 150 4. Final Dataset

151

152 The final dataset comprises 22,102 assemblies that belong to non-Typhoidal *Salmonella*. The  
153 final metadata table contains the following:

- 154 • Assembly ID: name of assembly ID as identified in the database;
- 155 • Database: name of repository from which said assembly was recovered;
- 156 • Collection Date: isolation date;
- 157 • Year: isolation year;
- 158 • ISO3: 3 letter code of the country of isolation;
- 159 • Latitude and Longitude: coordinates in decimal degree;
- 160 • Serovar: *Salmonella's* Serovar;
- 161 • ST: *Salmonella's* sequence type;
- 162 • BioSample: NCBI BioSample accession number;
- 163 • BioProject: NCBI BioProject accession number;
- 164 • Acquired Resistance: whether this assembly was found to contain ARGs or not;
- 165 • MDR Score: calculated MDR Score;

166

167

## 168 5. Model Weights Calculation

169 For the temporal trend analysis of resistance, we need to weight the observations relative to  
170 their representativeness in our dataset. To achieve this, we weighted all observations by the  
171 countries' Population Correction Unit (PCU) for each host (expressed as proportion) and  
172 corresponding isolation year times the proportion of genomes contributed by a given a  
173 country for a given year. We calculated PCU as described by Tiseo and colleagues (28) for all  
174 countries as follows:

$$175 \quad PCU_{k,s} = An_{k,s} \cdot (1 + n_{k,s}) \cdot \left( \frac{Y_k}{R_{CW/LW,k}} \right)$$

176 where  $An_{k,s}$  is the number of animal type,  $k$ , for each production system,  $s$  (intensive or  
177 extensive), in each country;  $n_{k,s}$  is the number of production cycles for each animal type in  
178 each production system;  $Y_k$  is the quantity of meat in each country for each animal type; and  
179  $R_{CW/LW,k}$  is the carcass weight to live weight ratio for each animal type. The PCU allows for  
180 direct comparisons of animals raised for food in across countries. For some countries, PCU  
181 data was unavailable before 1999 for Belgium, before 1991 for Belarus, before 1992 for Czech  
182 Republic and Slovakia before, and before 1991 for Belarus, Croatia, Estonia, and Lithuania. In  
183 addition, no PCU data existed prior to 1985. In such cases, we assigned the PCU value  
184 corresponding to the earliest available year in the time series. PCU data can be found in  
185 Supplementary Dataset S1.

186

187



188 **6. References**

189

190 1. Home - Nucleotide - NCBI.

191 2. Enterobase.

192 3. PATRIC.

193 4. Taxonomy - NCBI.

194 5. R Core Team (2019),. R: The R Project for Statistical Computing.

195 6. Wickham H, Averick M, Bryan J, Chang W, McGowan L, François R, Golemund G, Hayes

196 A, Henry L, Hester J, Kuhn M, Pedersen T, Miller E, Bache S, Müller K, Ooms J, Robinson

197 D, Seidel D, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H. 2019.

198 Welcome to the Tidyverse. J Open Source Softw 4:1686.

199 7. Dowle M, Srinivasan A, Gorecki J, Chirico M, Stetsenko P, Short T, Lianoglou S, Antonyan

200 E, Bonsch M, Parsonage H, Ritchie S, Ren K, Tan X, Saporta R, Seiskari O, Dong X, Lang

201 M, Iwasaki W, Wenchel S, Broman K, Schmidt T, Arenburg D, Smith E, Cocquemas F,

202 Gomez M, Chataignon P, Groves D, Possenriede D, Parages F, Toth D, Yaramaz-David

203 M, Perumal A, Sams J, Morgan M, Quinn M, @javrucebo, @marc-outins, Storey R,

204 Saraswat M, Jacob M, Schubmehl M, Vaughan D. 2019. data.table: Extension of

205 "data.frame."

206 8. Wickham H. 2019. plyr: Tools for Splitting, Applying and Combining Data.

207 9. Kahle D, Wickham H. 2013. ggmap: Spatial Visualization with ggplot2. R J 5:144.

- 208 10. Google Earth.
- 209 11. Spinu V, Grolemond G, Wickham H, Lyttle I, Costigan I, Law J, Mitarotonda D,  
210 Larmarange J, Boiser J, Lee CH. 2020. lubridate: Make Dealing with Dates a Little Easier.
- 211 12. Eddebuettel D. 2020. anytime: Anything to “POSIXct” or “Date” Converter.
- 212 13. O’Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, Rajput B, Robbertse  
213 B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V,  
214 Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D,  
215 Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey  
216 KM, Murphy MR, O’Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C,  
217 Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C,  
218 Webb D, Wu W, Landrum MJ, Kimchi A, Tatusova T, DiCuccio M, Kitts P, Murphy TD,  
219 Pruitt KD. 2016. Reference sequence (RefSeq) database at NCBI: current status,  
220 taxonomic expansion, and functional annotation. *Nucleic Acids Res* 44:D733–D745.
- 221 14. GitHub - PATRIC3/PATRIC-distribution: For distributing patric related code.
- 222 15. BLAST - Basic Local Alignment Search Tool. The Comprehensive Antibiotic Resistance  
223 Database.
- 224 16. BLAST: Basic Local Alignment Search Tool.
- 225 17. Bonomo RA. 2017.  $\beta$ -Lactamases: A Focus on Current Challenges. *Cold Spring Harb*  
226 *Perspect Med* 7.
- 227 18. Jacoby GA. 2009. AmpC  $\beta$ -Lactamases. *Clin Microbiol Rev* 22:161.

- 228 19. Evans BA, Amyes SGB. 2014. OXA  $\beta$ -Lactamases. Clin Microbiol Rev 27:241.
- 229 20. Salipante SJ, Hall BG. 2003. Determining the limits of the evolutionary potential of an  
230 antibiotic resistance gene. Mol Biol Evol 20:653–659.
- 231 21. Neuert S, Nair S, Day MR, Doumith M, Ashton PM, Mellor KC, Jenkins C, Hopkins KL,  
232 Woodford N, de Pinna E, Godbole G, Dallman TJ. 2018. Prediction of Phenotypic  
233 Antimicrobial Resistance Profiles From Whole Genome Sequences of Non-typhoidal  
234 *Salmonella enterica*. Front Microbiol 9.
- 235 22. Zankari E, Allesoe R, Joensen KG, Cavaco LM, Lund O, Aarestrup FM. 2017. PointFinder: a  
236 novel web tool for WGS-based detection of antimicrobial resistance associated with  
237 chromosomal point mutations in bacterial pathogens. 10. J Antimicrob Chemother  
238 72:2764–2768.
- 239 23. 2021. staramr. Python, National Microbiology Laboratory.
- 240 24. Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, Grundman H,  
241 Hasman H, Holden MTG, Hopkins KL, Iredell J, Kahlmeter G, Köser CU, MacGowan A,  
242 Mevius D, Mulvey M, Naas T, Peto T, Rolain J-M, Samuelsen  $\emptyset$ , Woodford N. 2017. The  
243 role of whole genome sequencing in antimicrobial susceptibility testing of bacteria:  
244 report from the EUCAST Subcommittee. Clin Microbiol Infect 23:2–22.
- 245 25. Kahlmeter G. 2015. The 2014 Garrod Lecture: EUCAST – are we heading towards  
246 international agreement? J Antimicrob Chemother 70:2427–2439.
- 247 26. Bush K. 2018. Past and Present Perspectives on  $\beta$ -Lactamases. 10. Antimicrob Agents  
248 Chemother 62:e01076-18.

- 249 27. Magiorakos A-P, Srinivasan A, Carey RB, Carmeli Y, Falagas ME, Giske CG, Harbarth S,  
250 Hindler JF, Kahlmeter G, Olsson-Liljequist B, Paterson DL, Rice LB, Stelling J, Struelens  
251 MJ, Vatopoulos A, Weber JT, Monnet DL. 2012. Multidrug-resistant, extensively drug-  
252 resistant and pandrug-resistant bacteria: an international expert proposal for interim  
253 standard definitions for acquired resistance. Clin Microbiol Infect Off Publ Eur Soc Clin  
254 Microbiol Infect Dis 18:268–281.
- 255 28. Tiseo K, Huber L, Gilbert M, Robinson TP, Van Boeckel TP. 2020. Global Trends in  
256 Antimicrobial Use in Food Animals from 2017 to 2030. 12. Antibiotics 9:918.
- 257 29. Kassambara A. 2021. rstatix: Pipe-Friendly Framework for Basic Statistical Tests.
- 258 30. Ebbert D. 2019. chisq.posthoc.test: A Post Hoc Analysis for Pearson’s Chi-Squared Test  
259 for Count Data.

260

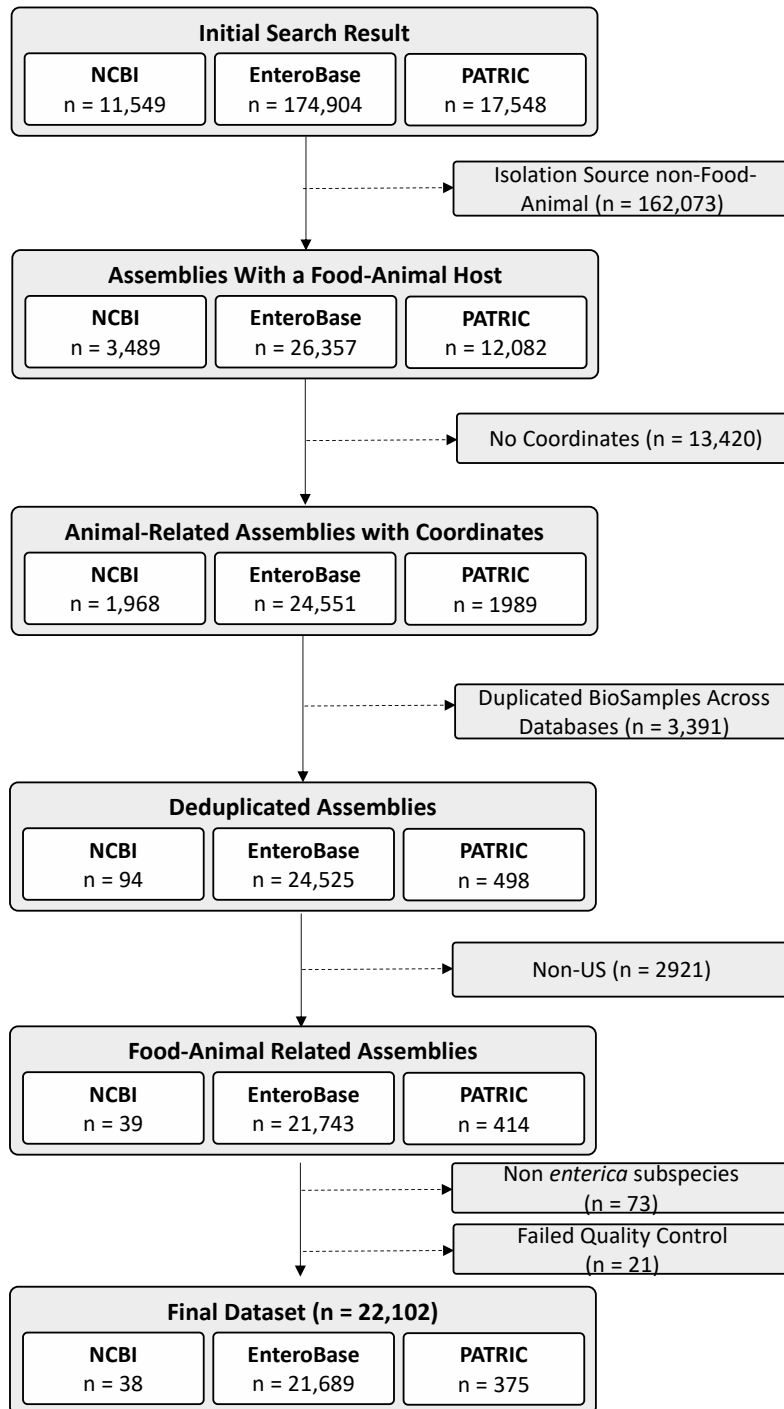
261

262

263 **II. Appendix Figures**

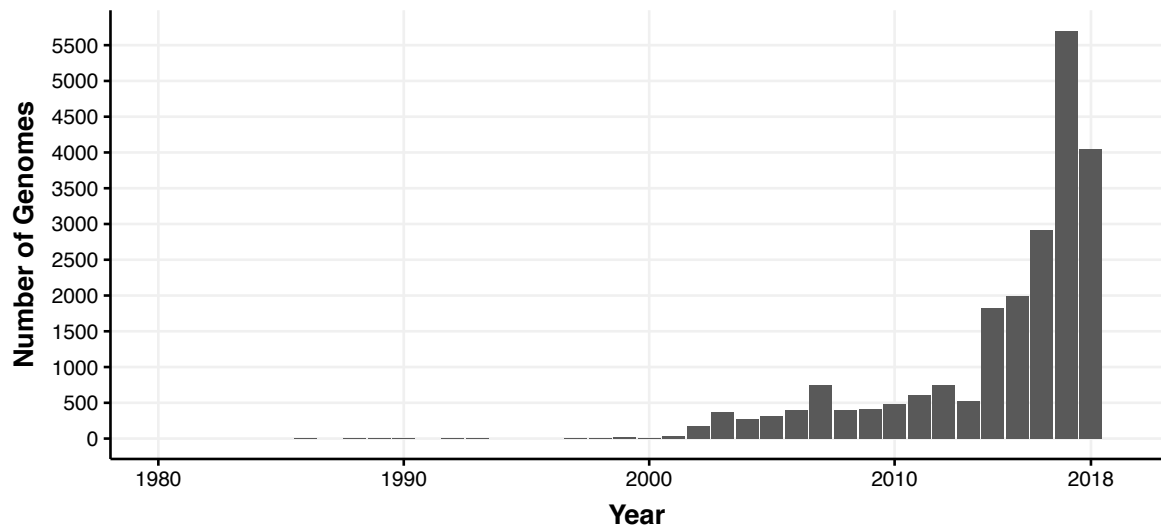
264

265



266  
 267 **Figure S1.** Number of genomes identified in public repositories and number of genomes  
 268 excluded throughout the curation process. NCBI - National Center for Biotechnology  
 269 Information (NCBI) Nucleotide database; PATRIC - Pathosystems Resource Integration  
 270 Center.

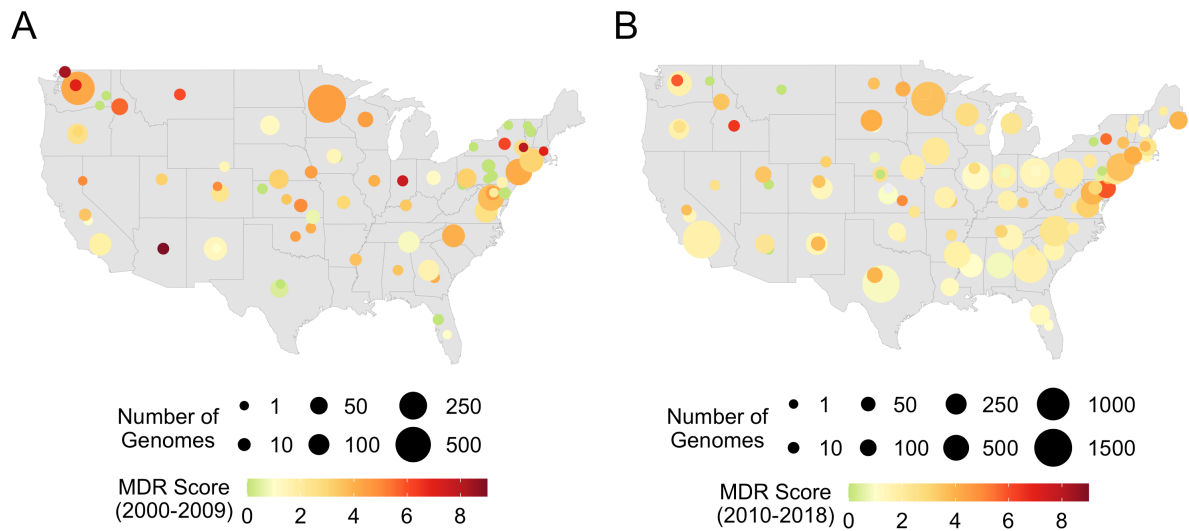
271  
 272  
 273



274

275 **Figure S2.** Distribution of the number of genomes per year.

276



277

278 **Figure S3.** Distribution of the of Multidrug Resistance Score (MDR Score) across the United

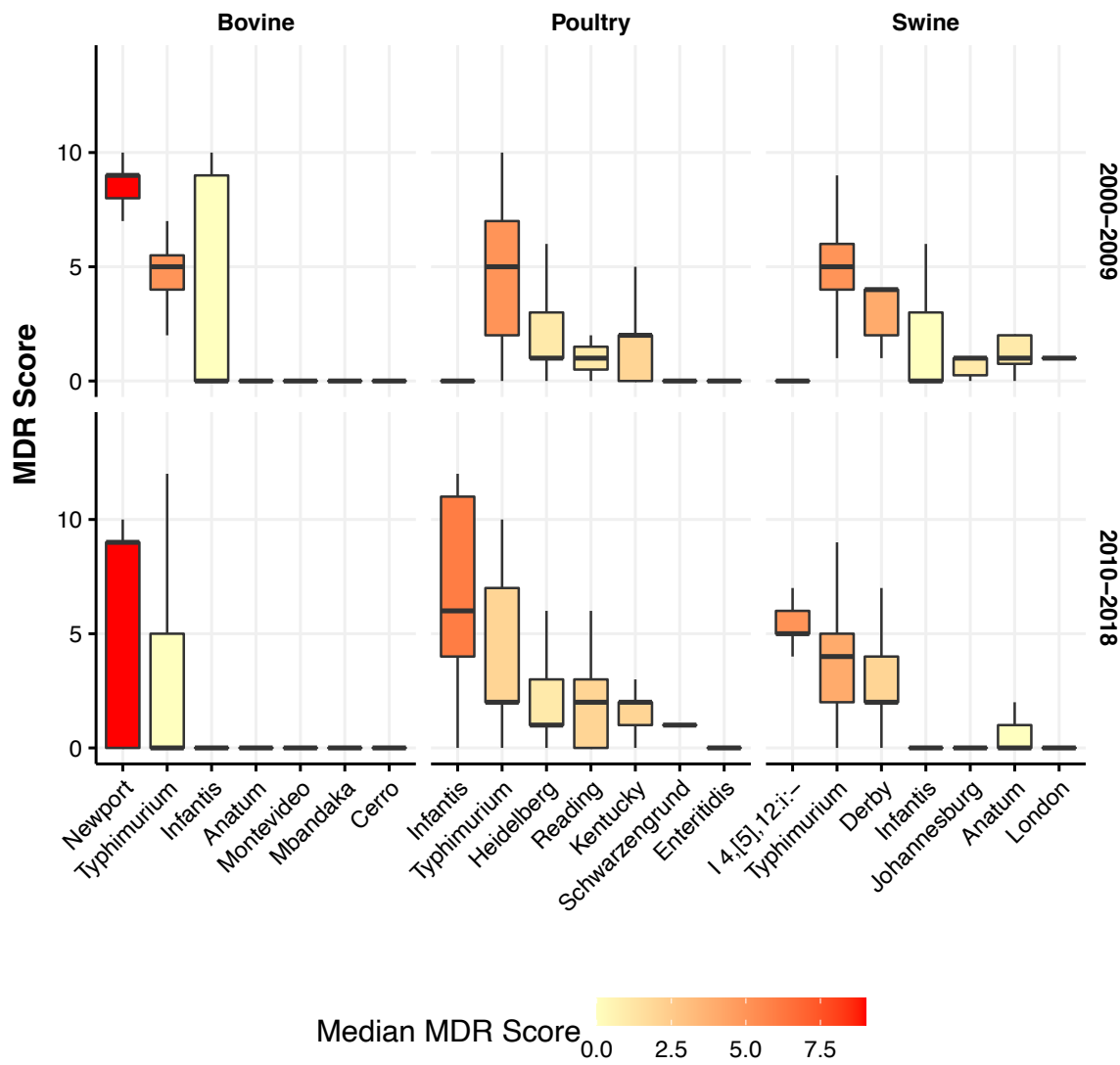
279 States. The dot size represents the number of genomes available for a single geographic

280 coordinate. **A.** Distribution of the MDR Score between 2000-2009; **B.** Distribution of the MDR

281 Score between 2010-2018.

282

283



284

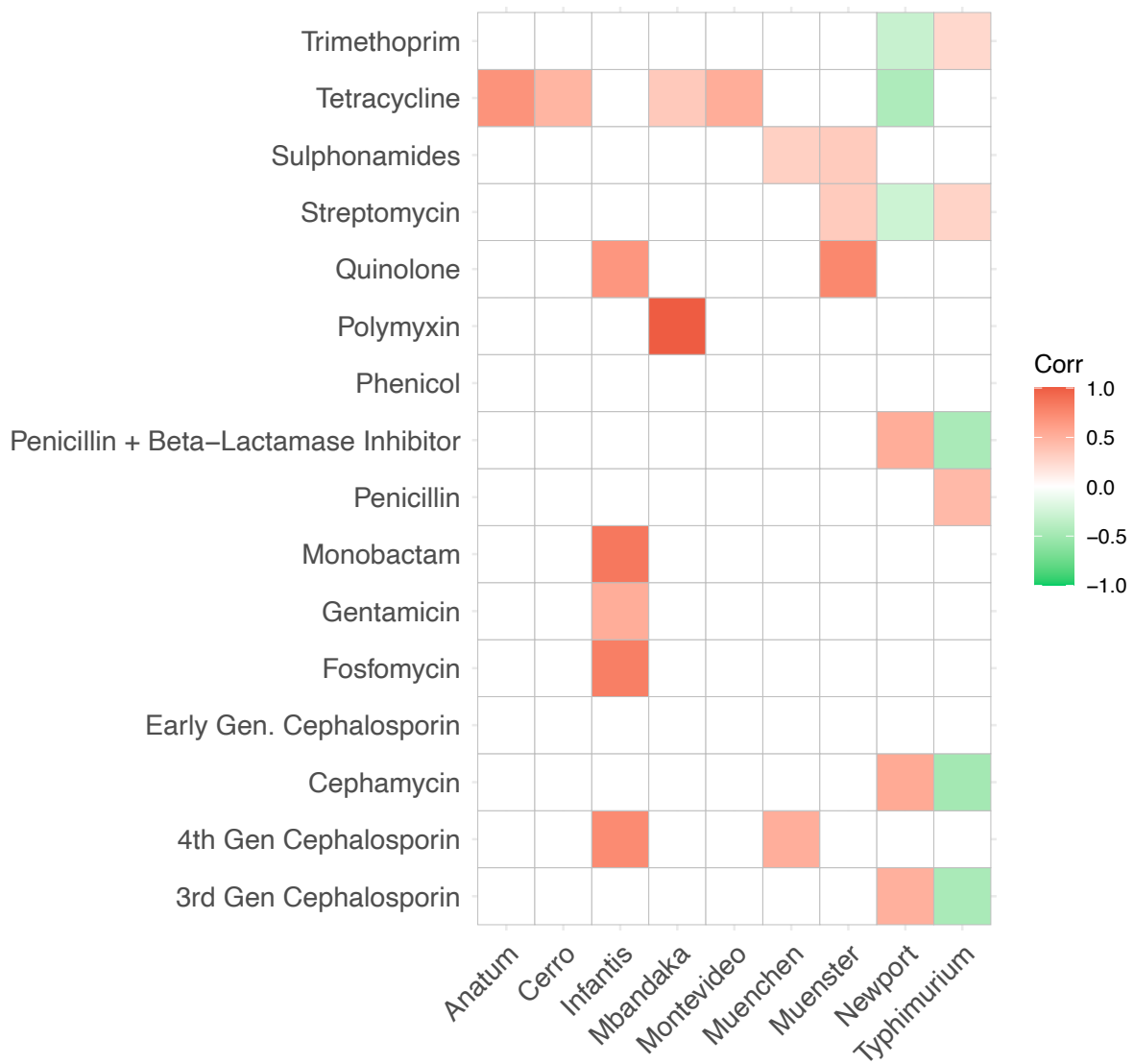
285 **Figure S4.** Distribution of the Multi Drug Resistance Score (MDR Score) per host per serovar

286 in 2000s (2000-2009) and 2010s (2010-2018).

287

288

289

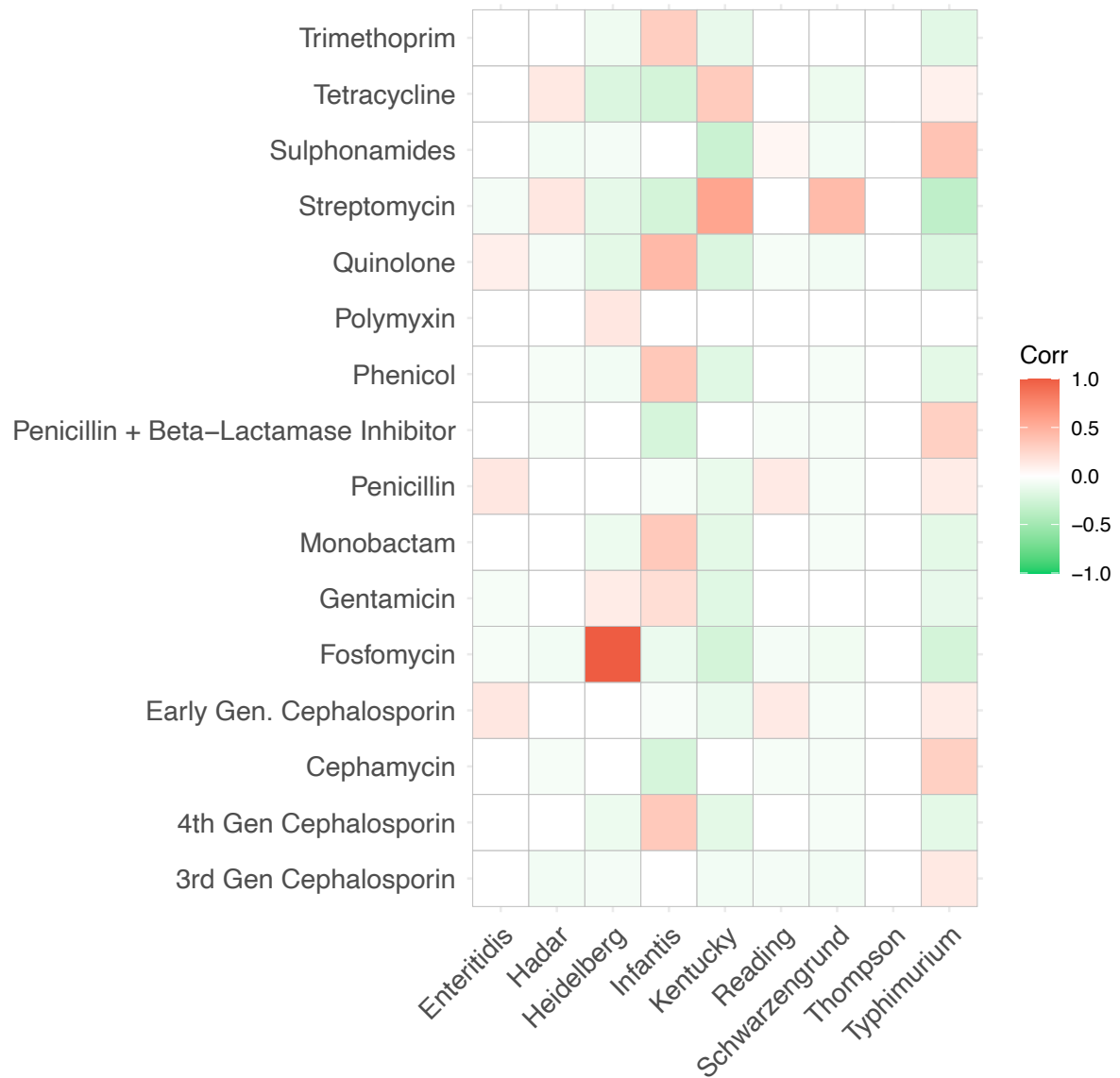


291

292 **Figure S5.** Correlation plot between the most frequent serovars in bovine and resistance293 phenotypes. Only correlations with an adjusted  $p$  value below 0.05 are shown. Non-

294 significant correlations are displayed as blank squares.





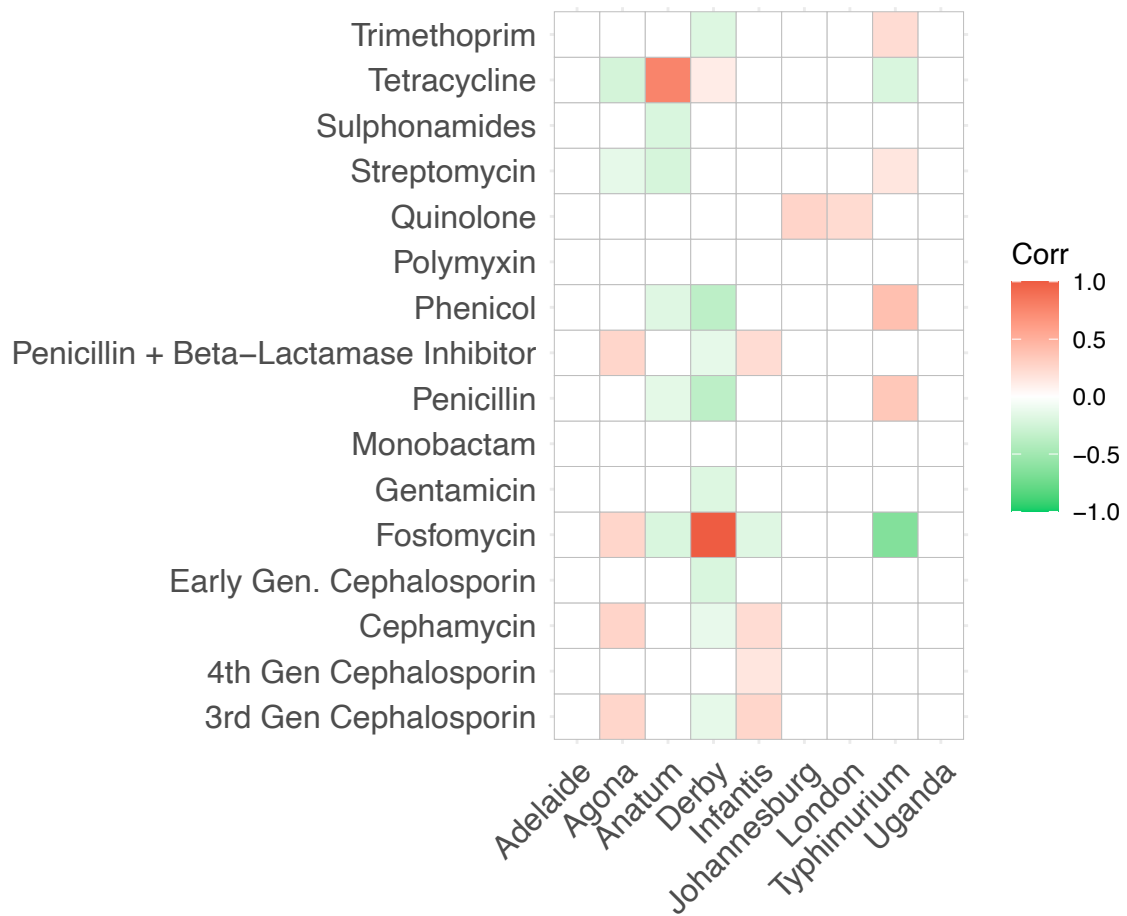
295

296 **Figure S6.** Correlation plot between the most frequent serovars in poultry and resistance

297 phenotypes. Only correlations with an adjusted  $p$  value below 0.05 are shown. Non-

298 significant correlations are displayed as blank squares.

299



300

301 **Figure S7.** Correlation plot between the most frequent serovars in swine and resistance

302 phenotypes. Only correlations with an adjusted *p* value below 0.05 are shown. Non-

303 significant correlations are displayed as blank squares.

304

305

306

307

308        **III.        Legends for Appendix Tables**

309

310

311        **Table S1.** Metadata file for the 22,102 genomes included for the dataset. Assembly\_ID –  
312 assembly identifier from the original database; Database – database from which assembly  
313 was retrieved; Collection\_Date – isolation date; Year – isolation year; Country - isolation  
314 country; ISO3 – three letter country code; Latitude – latitude geographic coordinate;  
315 Longitude – longitude geographic coordinate; Generic\_Host – food-animal host;  
316 Source\_Niche – indicates whether samples derive from food or from the animal itself;  
317 Source\_Type – specific animal species; Source\_Details – details as available in the database  
318 of origin; ST – *Salmonella* Sequence Type; Serovar – *Salmonella* Serovar; BioSample –  
319 National Center for Biotechnology Information BioSample accession number;  
320 Number\_Contigs – number of contigs in assembly; BioProject – National Center for  
321 Biotechnology Information BioProject accession number; MDR\_Score – calculated Multidrug  
322 Resistance Score; Acq\_Resist – whether assembly contains acquired resistance gene or not.

323

324        **Table S2.** Output from ResFinder. File\_Name – assembly name; Contig – contig name; Start –  
325 start position in the contig of the gene identified; End – end position in the contig of the gene  
326 identified; Gene – antimicrobial resistance gene identified; Coverage – proportion of gene  
327 present in the sequence; Coverage\_Map – visual representation of alignment of our sequence  
328 against the reference; Gaps – gaps in the sequence versus the reference; Perc\_Coverage –  
329 proportion of the gene covered; Perc\_Identity – proportion of nucleotide matches against  
330 reference; Database – reference database; Accession - National Center for Biotechnology  
331 Information accession number; Product – gene product; Class – predicted resistance to  
332 antimicrobial classes; Phenotype – predicted resistance to individual antimicrobials;

333 Mechanism of resistance – ResFinder specification of mechanism of resistance if available;  
334 Notes – further notes provided by the ResFinder on specific genes; Required\_gene – genes  
335 required to cause resistance phenotype if any. Gene\_clean – Harmonized gene name.

336 Table S2 can be found in the Zenodo repository in the following link:  
337 <https://zenodo.org/record/5519129#.YUzEj21Bw4g>

338

339

340 **Table S3.** Output from staramr PointFinder module. Assembly\_ID - assembly identifier from  
341 the original database, Gene – gene identified with mutation. Mutation designation in  
342 brackets; Type – mutation type; Position – amino acid position where mutation occurred;  
343 Mutation – specific mutation; Perc\_Identity – proportion of nucleotide matches against  
344 reference; Perc\_Overlap – proportion of the overlap between query and reference; HSP  
345 Length/Total Length – high scoring pair length over the length of the gene; Contig – contig  
346 name; Start – start position in contig; End – end position in contig.

347

348 **Table S4.** Fitted MDR Score values for all years and hosts. Year – isolation year; Generic\_host  
349 – animal host; mdr\_score – fitted MDR Score; se.fit - standard error; upp\_95 – upper bound  
350 of 95% confidence interval; low\_95 – lower bound of 95% confidence interval.

351

352 **Table S5.** Fitted antimicrobial resistance prevalence for individual classes for all years, hosts.  
353 Year - isolation year; Generic\_Host – animal host; Phenotype - antimicrobial class; Prevalence  
354 – fitted prevalence; low\_CI – lower bound of 95% confidence interval; upp\_CI – upper bound  
355 of 95% confidence interval. signif – whether covariate “Year” was statistically significant or not;

356

357 **Table S6.** Fitted antimicrobial resistance genes' prevalence for all years, hosts. Year – isolation  
358 year; Generic\_Host – animal host; Gene\_Dummy – acronym used to identify antimicrobial  
359 resistance gene; Prevalence – fitted prevalence; se.fit – standard error; Gene\_clean –  
360 antimicrobial resistance gene.

361

362 **Table S7.** Fitted serovar prevalence for all years, hosts. Year – isolation year; Generic\_Host –  
363 animal host; Serovar – *Salmonella* serovar; Prevalence – fitted prevalence; se.fit – standard  
364 error.

365

366 **Table S8.** Fitted serovar prevalence for 2018 and hosts. Year – isolation year; Generic\_Host –  
367 animal host; Serovar – *Salmonella* serovar; Prevalence – fitted prevalence; se.fit – standard  
368 error.

369

#### 370 **IV. Legend for Dataset S1**

371

372

373 **Dataset S1.** PCU data for all countries between 1985 and 2018. Each column corresponds to  
374 a food-animal/year combination. "Ca" refers to bovine, "Ch" refers to poultry, and "Pg" refers  
375 to swine.

376