

1 The genomic and epidemiological virulence patterns of *Salmonella*
2 *enterica* serovars in the United States

3

4 Gavin J. Fenske¹, Jane G. Pouzou¹, Régis Pouillot¹, Daniel D. Taylor¹, Solenne Costard¹, and
5 Francisco J. Zagmutt^{1*}

6

7 ¹EpiX Analytics, Fort Collins, Colorado, United States of America

8

9 *Corresponding Author

10 Email: fzagmutt@epixanalytics.com

11

12

13

14

15

16

17

18 **Abstract**

19 The serovars of *Salmonella enterica* display dramatic differences in pathogenesis and
20 host preferences. We developed a process (patent pending) for grouping *Salmonella* isolates and
21 serovars by their public health risk to provide better *Salmonella* control targets along the food
22 chain. We collated a curated set of 12,337 *S. enterica* isolate genomes from human, beef, and
23 bovine sources in the US. After annotating a virulence gene catalog for each isolate, we used
24 unsupervised random forest methods to estimate the proximity (similarity) between isolates
25 based upon the genomic presentation of putative virulence traits. We then grouped isolates
26 (virulence clusters) using hierarchical clustering (Ward's method), used non-parametric
27 bootstrapping to assess cluster stability, and externally validated the virulence clusters against
28 epidemiological virulence measures from FoodNet, the National Outbreak Reporting System
29 (NORS), and US federal sampling of beef products. We identified five stable virulence clusters
30 of *S. enterica* serovars. Cluster 1 (higher virulence) serovars yielded an annual incidence rate of
31 domestically acquired sporadic cases roughly one and a half times higher than the other four
32 clusters combined (Clusters 2-5, lower virulence). Compared to other clusters, cluster 1 also had
33 a higher proportion of infections leading to hospitalization and was implicated in more
34 foodborne and beef-associated outbreaks, despite being isolated at a similar frequency from beef
35 products as other clusters. We also identified subpopulations within 11 serovars. Remarkably, we
36 found *S. Infantis* and *S. Typhimurium* subpopulations that significantly differed in genome
37 length and clinical case presentation. Further, we found that the presence of the pESI plasmid
38 accounted for the genome length differences between the *S. Infantis* subpopulations. Our results
39 demonstrate that *S. enterica* strains with the highest incidence of human infections share a

40 common virulence repertoire. This work could be used in combination with foodborne
41 surveillance information to best target serovars of public health concern.

42

43 **Introduction**

44 Members of *Salmonella enterica* subspecies *enterica* are some of the most ubiquitous
45 agents implicated in foodborne human illnesses. Despite being constituents of the same
46 subspecies, members of *S. enterica* are not only commonly isolated from livestock but also
47 amphibians [1] and wild birds [2]. The wide host range for *S. enterica* makes control of the
48 pathogen exceedingly difficult due to the large number of potential reservoirs. Historically,
49 strains of *S. enterica* have been grouped into units termed serovars based upon serological
50 antigen presentation. While an initial list presented 44 *S. enterica* serovars in 1934 [3], today's
51 descriptions include over 2,500 serovars of *S. enterica* [4]. Nonetheless, in the US only 20
52 serovars accounted for 69.2% of human *S. enterica* isolates collected in 2016 by the US Centers
53 for Disease Control and Prevention's (CDC) Laboratory-based Enteric Disease Surveillance
54 (LEDS) program [5]. Furthermore, nearly 10% of *S. enterica*'s serovars may be polyphyletic or
55 paraphyletic [6].

56 To establish infections in disparate hosts, *S. enterica* manipulates common immune
57 functions of higher vertebrates. Indeed, the classic gastroenteritis associated with *S. enterica*
58 infections is the result of the pathogen affecting the host's innate immune system to generate
59 inflammation, subsequently producing unique metabolic niches for *S. enterica* while killing its
60 competitors for reduced substrates in the hindgut [7-9]. Such remarkable expropriation of the
61 hosts immune functions is achieved by virulence genes (virulence factors), many of which are
62 contained within chromosomal elements termed Salmonella Pathogenicity Islands (SPI) [10].

63 Genes contained within SPI aid in host cell invasion, and subsequent survival and dissemination
64 within and between Eukaryotic host cells [11,12]. However, serovars display differences in
65 pathogenesis and host-preferences. For example, the human-restricted serovar *S. enterica* ser.
66 Typhi (*S. Typhi*), the etiological agent of typhoid fever, does not typically cause submucosal
67 inflammation and resultant diarrhea in infected patients as with classical salmonellosis, but
68 instead elicits a systemic enteric fever characterized by initial immune evasion [13,14]. *S.*
69 Dublin, a bovine adapted serovar, commonly generates systemic infections in humans and is
70 isolated from blood samples in 61% of human clinical infections as compared to an average of
71 5% for other *S. enterica* serovars in the US [15]. The general pathogenesis of *S. enterica* is not
72 fully elucidated, and the virulence potential for individual serovars is poorly understood.
73 Furthermore, most studies have focused upon *S. Typhimurium* as a model organism for all *S.*
74 *enterica* virulence, [8,16-21] which could obfuscate differences between serovars.

75 Despite the tremendous virulence diversity within *S. enterica*, microbial criteria from the
76 US Food Safety and Inspection Services (FSIS) on important sources of *S. enterica* such as beef
77 and poultry meats target all serovars equally, based on prevalence. Further, traditional
78 surveillance methods can take considerable time to identify emerging serovars of public health
79 concern, thereby delaying food safety intervention implementation [22]. Understanding virulence
80 differences between serovars and identifying emerging virulent serovars in a timely manner can
81 inform more focused risk management strategies targeting serovars with an inordinate impact on
82 public health while reducing food waste due to recalls.

83 Previous studies have used genomics to identify serovar groups of public health concern.
84 Karanth et al. analyzed a limited number of genomes and serovars originating from humans,
85 poultry, and swine to characterize virulent serovars [23]. This analysis had the benefit of using

86 the entire genome of *Salmonella* to group isolates by disease presentation; however, the
87 computational resources required prevent its application to a large number of isolates. In another
88 study, researchers used single nucleotide polymorphism (SNP) clusters and *S. Saintpaul* as a
89 model to identify virulent isolates [24]. Although using high-resolution genomic methods
90 identified SNP clusters associated with a high proportion of human clinical isolates, *S. Saintpaul*
91 may not be the best serovar model due to its polyphyletic nature [25]. The objective of this study
92 was to develop a computationally efficient genomic approach to group *Salmonella* serovars by
93 virulence biomarkers in isolates from humans, beef, and bovine animals and define the human
94 health risk of the resulting clusters using epidemiological data. We chose beef as a model
95 foodstuff since US federal monitoring of *Salmonella* in beef is well-established, nationally
96 representative, and beef remains an important vehicle for *S. enterica*. Beef production in the US
97 is more decentralized than poultry and pork production [26] and we hypothesize that this
98 decentralization may present unique genomic populations arising from geographic separation as
99 previously observed in *S. enterica* serovars [27,28]. Furthermore, *S. enterica* in beef products is
100 understudied compared to other vehicles such as eggs, poultry, and pork meat.

101

102 **Materials and methods**

103 We developed an information pipeline (patent pending) using virulence factors as
104 markers and epidemiological data as validation to group serovars by their risk to human health .
105 After compiling a curated set of *S. enterica* genomes (n=12,337) from human, bovine, and beef
106 sources, we applied an unsupervised random forest and hierarchical clustering approach to group
107 isolates based upon genomic virulence trait presentation and validated the groups against

108 epidemiological measures including clinical presentation from sequenced isolates collected by
109 the FoodNet active surveillance network (29).

110 **Contig assembly selection and quality criteria**

111 We compiled *S. enterica* assemblies from bovine-associated isolates from three primary
112 sources: 1. BioProject PRJNA242847 (FSIS HACCP samples, accessed 7/13/2021), 2.
113 BioProject PRJNA292666 (FSIS NARMS isolates, accessed 7/13/2021), and 3. BioProject
114 PRJNA292661 (FDA NARMS isolates, accessed 8/25/2021). We collected isolates from sources
115 specified as bovine-associated or beef origin from the metadata for the above BioProjects.

116 We retrieved *S. enterica* isolates from human clinical cases from BioProject
117 PRJNA230403 (CDC PulseNet, accessed 9/13/2021) and identified sporadic, domestically
118 acquired, non-attributed *S. enterica* isolates from the FoodNet active surveillance network [29].
119 We did not include outbreak cases from FoodNet since they are not attributed to a particular
120 source in that dataset. Instead, we used the National Outbreak Reporting System (NORS) dataset,
121 a passive system for reporting enteric disease outbreaks in the US, to identify additional human
122 isolates that originated from beef-attributed outbreaks. We initially defined beef attribution based
123 on the Interagency Food Safety Analytics Collaboration (IFSAC) classification. As IFSAC-
124 defined beef-associated salmonellosis outbreaks for which clinical isolates were sequenced are
125 limited, we widened the definition of potentially beef-associated illness to include outbreaks
126 which listed beef as an identified contaminated ingredient. We based this inclusion on whether
127 the list of foods or ingredients per outbreak included beef, even if other possible ingredients
128 could not be ruled out to definitively assign an IFSAC classification. The following text strings
129 were used to identify beef-associated outbreaks: "beef", "burger", "steak", "carne", "kitfo", "ox

130 tongue", "short-rib", "prime rib", "barbacoa". If the IFSAC classification attributed an outbreak
131 to other foods, we did not designate it as a beef-associated outbreak.

132 We removed isolates from the data set if: 1) No pre-computed assembly was available on
133 NCBI, 2) SKESA v. 2.2 was not used to construct the assembly, 3) The number of contigs
134 representing the assembly was greater than 300, and 4) The contig n50 was less than 25,000 bp.
135 After initial parsing for isolation sources and assembly quality, we included serovars with 50 or
136 more isolates in the analysis. In total, the final analysis set includes 12,337 assemblies and
137 represents 37 serovars.

138 **In silico serovar prediction**

139 We used *Salmonella in silico* Typing Resource (SISTR) [30] with default options to
140 assign putative serovars to each assembly. 1,077 assemblies failed the quality control step within
141 the SISTR software with the same error message "FAIL: Wzx/Wzy genes missing...", but all
142 330 / 330 genes for the core genome multilocus sequence typing (cgMLST) scheme used within
143 the software were present within these assemblies. We retained assemblies failing QC with the
144 aforementioned error message and that contained all 330 cgMLST loci for the analysis while
145 excluding any assemblies which failed the quality control step and did not have all cgMLST
146 genes. S1 table provides the SISTR serotyping results for each assembly.

147 **Virulence gene annotation**

148 We collated a custom database of putative virulence factors associated with *Salmonella*,
149 *Escherichia*, *Shigella*, and *Yersinia* from the virulence factor database (VFDB) (accessed
150 9/13/2021) [31] and putative virulence factors associated with *Salmonella*, *Escherichia*, and
151 *Shigella* from PATRIC (accessed 9/13/2021) [32]. Next, we combined amino acid sequences of
152 the open reading frames (ORF) with a reference proteome of *Salmonella* Typhimurium LT2

153 (<https://www.uniprot.org/proteomes/UP000001014>) and made the database non-redundant by
154 clustering the open reading frames at 0.90 global identity using cd-hit [33]. We passed the
155 resultant database to Prokka [34] using the "--proteins" option to specify the database as the
156 primary annotation database in the software pipeline. We then parsed gene annotations from the
157 VFDB and PATRIC non-redundant database from the resultant Prokka annotation tables.
158 Additionally, to ensure consistent ORF predictions between assemblies, we trained a model
159 using Prodigal [35] on the chromosome of the reference Salmonella Typhimurium LT2 assembly
160 ASM694v2 (https://www.ncbi.nlm.nih.gov/assembly/GCF_000006945.2/) and passed the
161 training file to Prokka using the command "--prodigaltf".

162 **pESI plasmid identification**

163 The pESI plasmid is an emerging concern in some *S. enterica* serovars, namely *S.*
164 *Infantis* [36]. To confirm the presence of this plasmid, we compiled an additional nucleotide
165 database of 13 marker genes previously used to identify pESI plasmids in *S. enterica* contig
166 assemblies from two sources [36,37]. We then conducted a nucleotide BLAST search against the
167 database using the software Abricate [38] and defined positive hits as a percent identify of $\geq 95\%$
168 and percent coverage of $\geq 95\%$. We presumed that an isolate contained the pESI plasmid if the
169 contig assembly was positive for the pESI specific *repA* gene and contained at least five
170 additional marker genes.

171 **Random forest model construction**

172 We used virulence factor gene annotations from the resultant Prokka outputs and
173 constructed a count matrix of virulence genes (assemblies x virulence factors). We excluded
174 putative virulence factors present in more than 95% of assemblies or which were not present in at
175 least 10 (0.08%) assemblies. To generate row similarity (isolate relatedness), we fit an

176 unsupervised random forest to the count matrix of virulence factors using the randomForest [39]
177 package in R[40]. The random forest model contained 50,000 trees and used 60 features
178 (columns, virulence factors) at every split.

179 **Grouping isolates and assessing cluster stability**

180 We converted the row-wise proximity matrix from the random forest model to a distance
181 matrix (1 – similarity) and subjected it to agglomerative clustering using Ward’s method [41].
182 We conducted two analyses: 1) $k = 5$, used to group serovars by putative virulence factor
183 patterns and 2) $k = 37$, using the same number of clusters as serotypes in the data to identify
184 possible subpopulations within serovars with different virulence factor catalogues. We defined
185 cluster stability for both analyses as a Jaccard similarity of ≥ 0.75 [42] for 10,000 non-parametric
186 bootstrap samples. For the main serovar clustering analysis, we choose $k=5$ as it was the highest
187 value for which all clusters were stable. The $k=37$ version of the analysis tested whether
188 distribution of virulence factor combinations is more similar within each serotype than between
189 serotypes and, therefore, whether serotype is a reasonable representation of these virulence
190 differences. We defined serovar subpopulations as serovars with at least two of these populations
191 annotated into different clusters ($k = 37$ clustering) and with at least two of the populations
192 representing greater than or equal to 0.20 of the total serovar population. Clustering and
193 bootstrapping were performed using the clusterboot function from the fpc [43] package in R.

194 **Epidemiological indicators**

195 We estimated epidemiological indicators for both virulence clusters and serovars using
196 sporadic and domestically acquired cases from FoodNet (2016 – 2019). We excluded outbreak-
197 associated cases to decrease bias due to outbreak size and removed travel-related cases to
198 exclude foodstuffs from regions that may have different *S. enterica* population structures than the

199 US and are not good indicators of US consumer exposure to *S. enterica* in food. We calculated
 200 the proportion of positive samples for *Salmonella*, individual serovars, and each cluster in beef
 201 products from FSIS regulatory testing programs: MT43 (raw ground beef), MT60
 202 (manufacturing trimmings) and MT64 (Components other than trim) for 2016 – 2019. We
 203 modeled all proportions using $Beta(s+0.5, n-s+0.5)$ (eqn. 1), with a Jeffery’s prior ($Beta(0.5, 0.5)$)
 204 as a Bayesian conjugate to the Binomial distribution[44]. Table 1 lists parameters s and n used to
 205 model these proportions. We used 1M Monte Carlo simulations to estimate and compare
 206 posterior distributions using numerical integration, with a 99% confidence level for statistical
 207 significance.

208

209 **Table 1. Description of parameters used for calculation of epidemiological indicator**
 210 **estimations using Equation 1**

Model variable / Epidemiological indicator	Parameter description	
	s (numerator)	n (denominator)
Salmonella proportion positive in beef (p_+)	(s_+) : number of FSIS samples positive for Salmonella	(n_+) : number of FSIS samples taken
Hospitalization proportion	H_k : number of hospitalizations from cluster k ($k=1\dots5$)	I_k : number of total infections in FoodNet from cluster k ($k=1\dots5$)
Extraintestinal Infection proportion	EI_k : number of extraintestinal infection from cluster k ($k=1\dots5$)	
Mortality Proportion	M_k : number of deaths from cluster k ($k=1\dots5$)	
Proportion of gene presence within respective cluster group ($p_{g,k}$)	$n_{k,g}$: number of isolates from cluster (k =Cluster 1 vs. clusters 2-5) with gene g ($g=1\dots182$)	n_k number of isolates from cluster k (k =Cluster 1 vs. clusters 2-5)

211

212 Incidence of domestically acquired sporadic cases

213 We modeled the incidence of domestically acquired sporadic cases per 100k people per
 214 year (λ_{ij}) for serovar j in FoodNet state i using the Bayesian conjugate for a Poisson rate with
 215 Jeffrey’s prior $Gamma(0.5, 0.00001)$ (44), hence $Gamma(\alpha_i + 0.5, \beta_i t + 0.00001)$ (eqn. 2),
 216 where α_{ij} is the serovar case totals per state and β_i is the FoodNet catchment area population for

217 $t=4$ years. Sporadic cases are defined as illnesses which were not linked to a known outbreak.
218 The catchment area of FoodNet is not evenly distributed between the 10 participating states, so
219 the population-weighted mean serovar incidence (λ_j) for the study period was $\sum_{i=1}^{10} \lambda_{ij} p_i$, where
220 λ_{ij} is the FoodNet state-specific mean serovar incidence and p_i is the state catchment proportion
221 of total FoodNet population. Virulence cluster population-weighted average incidence, λ_k , is the
222 sum of the clusters' constituent serovars incidence rates for the study period ($\lambda_k = \sum_{j=1}^{j_k} \lambda_j$).

223 **Serovar proportion positive in beef**

224 We determined the proportion of *Salmonella* positives (p_+) following eqn. 1, with s_+
225 number of samples positive for *Salmonella* and n_+ total number of samples from FSIS testing.
226 We estimated the proportion of serovar j isolated from beef products (X_j) with a Dirichlet
227 distribution, ($Dir(a_j)$ eqn. 3), where a_j is the number of FSIS isolates from serovar j . We excluded
228 serovars without any positive isolates in FSIS testing and retained all serovars from the testing
229 program (including those not included in the analysis set).

230 Serovar proportion positive (p_{j+}) was taken as the product of total *Salmonella* proportion
231 positive (p_+) and the serovar proportion of the total *Salmonella* population from eqn. 3 (X_j).

232 Finally, we derived the cluster proportion positive (p_{k+}) as the summation of the cluster's
233 constituent serovars.

234 **Hospitalization, Extraintestinal Infection, and Mortality Proportions**

235 We determined the proportion of infections with a certain outcome (i.e., hospitalization,
236 extraintestinal infections, and mortality) for each cluster k ($k=1 \dots 5$) and for cluster 1 vs the
237 combinations of others. We defined extraintestinal infections as having "URINE", "BLOOD",

238 "ORTHO", "ABSCESS", "OTHER STERILE SITE" and "CSF" isolation sources. We modeled
239 all the proportions using eqn. 1, with parameters described in Table 1.

240 **Differential gene carriage**

241 To identify virulence factors differential between clusters, we trained a supervised
242 random forest model (ntree = 5,000, features to try at split = 13) to classify isolates into two
243 groups: cluster 1 (higher virulence) and clusters 2-5 (lower virulence). We extracted variable
244 importance from the random forest model and defined factor importance using the mean
245 decrease of Gini impurity. As with other proportions, we used eqn. 1 to model the proportion of
246 factor presence ($p_{g,k}$) within respective cluster group (1 vs. 2-5) for each of the virulence factors
247 used in the random forest. The relative frequency (RF) of a given factor was the resultant ratio of
248 proportions of factor presence ($RF = p_{g_1} / p_{g_{(2-5)}}$). To ease interpretation of differential genes, we
249 categorized them into five broad virulence mechanism categories (Adhesion, Motility/Invasion,
250 Survival/Host Persistence, Toxin and Virulence Factor Secretion) (Table S3) based on their
251 virulence descriptions listed in the VFDB [31] and PATRIC [32] databases.

252 **Code and data availability**

253 Aggregated data and code used to generate figures for this study are available in our
254 online repository: (link will be provided upon publication, Files shared with reviewers). FoodNet
255 data is used with permission from the Centers for Disease Control and Prevention and although
256 raw data may not be shared, code written from aggregated inputs is provided in our online
257 repository.

258

259 **Results**

260 **Genome assemblies analyzed**

261 The Pathogen Detection Network hosted by NCBI contains over 400,000 sequenced
262 *Salmonella* isolates from various sources and contributors. From these, we extracted 53,849
263 isolates from specific sampling programs. We further reduced to a final analysis set of 12,337 *S.*
264 *enterica* assemblies comprised of 37 serovars representing human clinical cases in the US and
265 bovine and beef associated isolates, (Fig 1). Approximately 55% (6,751) assemblies are from US
266 human clinical infections with the remaining 45% (5,586) representing isolates from bovine
267 animals and beef products. The metadata for the genomes analyzed is provided in S1 Table.

268 **Fig 1. Analysis set of genomes.** Description of the *S. enterica* genome assemblies considered,
269 and exclusion and inclusion criteria applied to generate the analysis set.

270 **Clustering serovars using isolate virulence gene catalogues**

271 To establish clusters of serovars, we identified virulence factors genes from each
272 assembly and compiled them into a count matrix, trained an unsupervised random forest model
273 to approximate similarity between isolate virulence factor catalogues (Fig 2), and subjected the
274 resultant isolate similarity matrix to agglomerative clustering to identify clusters with subsequent
275 non-parametric bootstrapping to validate cluster stability.

276 We identified five stable clusters of *S. enterica* isolates (Fig 3A), with the majority of
277 serovar isolates residing within the same clusters (mean within serovar cluster proportion =
278 0.96). (Fig 3B). However, *S. Reading* (cluster 1: n = 28 (0.47), cluster 3: n = 32 (0.53)), *S.*
279 *Saintpaul* (cluster 3: n = 134 (0.66), cluster 1: n = 68 (0.34)), and *S. I 1,4,[5],12:b:-* (cluster 1: n =
280 53 (0.66), cluster 3: n = 30 (0.34)) had at least 33% of total serovar isolates in two different

281 virulence clusters. The five virulence clusters are of uneven size (Fig 3C) with cluster 1
282 containing almost 10 times more assemblies than cluster 2. We attempted to decrease the size of
283 cluster 1 by introducing a sixth cluster. However, the sixth cluster was unstable (bootstrap
284 Jaccard similarity = 0.515) and cluster 4 was split, not cluster 1, indicating that the variance
285 (Ward's method used to cluster) within cluster 1 is less than that of cluster 4 despite its much
286 larger size (S1 Fig). Interestingly, cluster 2 is comprised of only *S. Javiana* and the cluster
287 homogeneity of *S. Javiana* was preserved with the addition of the sixth cluster (S1 Fig). Serovar
288 cluster designations are provided in Table S2.

289 **Fig 2. Conceptual model of virulence cluster development.** First, we downloaded contig
290 assemblies and quality controlled for fragmentation followed by the identification of virulence
291 factors. We then fit an unsupervised random forest model to the isolate level virulence factors
292 catalogues to approximate relatedness. We converted the resultant similarity matrix to a distance
293 matrix (1 – similarity) and clustered using Ward's method. We identified five stable clusters and
294 validated using non-parametric bootstrapping.

295 **Fig 3. Description of the five virulence clusters.** (A) Dendrogram depicting the hierarchical
296 relationship between 12,337 *S. enterica* genome assemblies based upon virulence factor gene
297 carriage with the five virulence clusters superimposed on top. (B) Heatmap of serovar proportion
298 within each of the five respective virulence clusters. Rows are clustered using Ward's method.
299 (C) Characteristics of the five virulence clusters: cluster stability - Jaccard similarity of 10,000
300 non-parametric bootstraps, Number of Genomes - depicting the number of *S. enterica* genomes
301 constituent in each cluster, and number of serovars (within cluster serovar proportion > 0.5) in
302 each cluster.

303 **General epidemiological characteristics of virulence clusters**

304 To investigate if the genomic virulence clusters correspond to clinical case presentation,
305 we computed basic epidemiological characteristics per cluster for 2016-2019 as proxies for
306 virulence phenotypes: proportion positive in beef products, number of outbreaks, incidence of
307 domestically acquired sporadic cases per 100k people per year, hospitalization proportion given
308 infection, extraintestinal infection proportion given infection, and mortality proportion given
309 infection. We computed the results by virulence cluster (S4 Table). and by serovar. Not every *S.*
310 *enterica* captured during surveillance programs in the US is subjected to sequencing, therefore
311 we attributed cases from a given serovar to the cluster to which the highest proportion of serovar
312 isolate was assigned (e.g., 98.5% of *S. Typhimurium* isolates were resident in cluster 1, therefore
313 all cases of *S. Typhimurium* in the datasets were allocated to cluster 1). Cluster 1 serovars have
314 the highest incidence rate of domestically-acquired sporadic cases (5.9 cases per 100k population
315 per year, 99% CrI: 5.77 – 6.06) (Fig 4A), approximately 1.5x higher than that of clusters 2-5
316 combined during 2016 – 2019 (incidence rate ratio: 1.5, 99% CrI: 1.44 – 1.55). Moreover,
317 infections from serovars in cluster 1 had a higher proportion of hospitalizations than serovars in
318 cluster 2 (relative frequency (RF): 1.10, 99% CrI: 1.002 – 1.200), cluster 4 (RF: 1.15, 99% CrI:
319 1.029 – 1.296), and cluster 5 (RF: 1.17, 99% CrI: 1.058 – 1.288) (Fig 4B). The cluster 1
320 proportion positive in beef products was less than half of clusters 3-5 (proportion positive ratio:
321 0.44, 99% CrI: 0.366 – 0.528) (Fig 4C). However, cluster 1 serovars were implicated in the
322 highest proportion of total foodborne outbreaks and beef associated outbreaks in the US from
323 2016 – 2019 (Fig 4D), generating approximately 2.5x more beef associated outbreaks (20 vs. 8)
324 than clusters 3-5 combined (There were no cluster 2 (i.e., *S. Javiana*) isolates found in beef
325 sampling or in beef associated outbreaks). Additionally, higher virulence serovars were involved

326 in approximately 1.47x more foodborne outbreaks than clusters 2-5 combined from 2016 - 2019
327 (285 vs. 194).

328

329 **Fig 4. Epidemiological indicators of the five virulence clusters for the study period 2016-**

330 **2019.** (A) Incidence of domestically acquired sporadic cases per 100k population per year by

331 virulence cluster. (B) Proportion of clinical infections resulting in hospitalization by virulence

332 cluster. (C) Proportion positive estimates in FSIS testing of US beef products. No isolates from

333 cluster 2 (comprised solely of *S. Javiana*) were retrieved from 2016 – 2019. (D) Proportion of

334 total US foodborne and beef-associated outbreaks attributed to serovars in the analysis set.

335

336 **Differential carriage of virulence factors between clusters**

337 In addition to the clear difference in epidemiological characteristics, a clear bifurcation

338 exists between cluster 1 and clusters 2-5 (Fig 3A). We sought to identify virulence mechanism

339 categories driving this differentiation by visually exploring the abundance (number of isolates

340 carrying at least one copy of the virulence gene), frequency (RF), and clustering influence (mean

341 Gini impurity) of the 182 virulence genes used in our analysis (Fig 5). The top two quadrants of

342 the figure include the virulence factors that provided the most separation between clusters

343 (highest mean Gini impurity), and the virulence factors on the two right quadrants were more

344 common in the higher virulence group than in the lower virulence group (higher RF).

345 Consequently, the upper-right quadrant includes factors that best distinguish the higher virulence

346 cluster and were most frequent in the higher virulence group. Interestingly, factors involved in

347 adhesion were the most important differentiators between clusters, while also being present in

348 both virulence groups. Expectedly, more abundant genes generally provided higher

349 differentiation, as some genes were very rare (e.g., 11 isolates contained the ompF gene), while
350 others, such as the ratB gene, were found in up to 95% of isolates, as dictated by our gene
351 exclusion criteria. Besides adhesion, other influential genes were involved in survival/host
352 persistence, and motility/invasion, whereas genes manifesting toxin production provided less
353 differentiation between virulence groups (Fig 5). The full list of virulence factors considered in
354 the analysis and gene metadata are provided in Table S3, while genes present in each isolate
355 (n=68) are presented in Table S5.

356

357 **Fig 5. Abundance, relative frequency, and influence of 182 virulence factors used to classify**
358 **12,337 *S. enterica* genomes from human, beef, and bovine sources into two virulence**
359 **groups.** Points represent the square root mean decrease Gini impurity and natural log relative
360 frequency (Cluster 1/Cluster 2-5) of each virulence factor, with diameter proportional to the
361 number of isolates (min: 11, max: 11,653) carrying at least one copy of the virulence gene. Color
362 designates virulence mechanism categories, as derived from the VFDB [31] and PATRIC [32]
363 databases. Vertical dashed line represents equal frequency of virulence factors between the
364 clusters, while points to the right of this line represent factors more frequently found in cluster 1.
365 Points above the horizontal dashed line (square root of the mean decrease Gini impurity=12.5)
366 represent virulence factors that were more influential differentiators in classifying isolates.

367

368 **Within serovar virulence subpopulations**

369 Horizontal gene transfer molds virulence gene carriage, especially within SPI [45,46].
370 We hypothesized that horizontal gene transfer may lead to virulence subpopulations that could
371 be identified using random forest methods otherwise missed in more traditional alignment-based

372 phylogeny methods. To test this hypothesis, we increased the number of clusters to correspond to
373 the number of serovars ($k = 37$). If no virulence subpopulations are present (within serovar
374 variance is less than between serovar variance), each of the 37 clusters should contain a majority
375 of one serovar (see methods). However, we found 11 serovars with virulence subpopulations
376 (Table 2). The full list of subpopulation designations is provided in S6 Table. To test if virulence
377 subpopulations may correspond to phenotypic differences in case presentation, we computed the
378 proportion of clinical infections resulting in extraintestinal infections for each serovar
379 subpopulation for sequenced strains with case presentation in the FoodNet surveillance system.
380 Two serovars yielded significant differences in invasiveness between serovar subpopulations. *S.*
381 *Infantis* split into two subpopulations (subpopulation 18: $n = 145$, subpopulation 20: $n = 243$) as
382 shown in Fig 6A. The genome assembly size for subpopulation 18 isolates was significantly
383 longer (4.98 Mb vs. 4.68 Mb, p -value $< 2.2E-16$, Mann-Whitney U test) (Fig 6B) than isolates
384 from subpopulation 20. Of the 388 *S. Infantis* genome assemblies in the analysis set, 242 had
385 associated clinical presentation data from FoodNet split evenly between the two subpopulations
386 ($n = 121$, $n = 121$). Isolates from subpopulation 18 were more than twice as likely to result in
387 extraintestinal clinical infections than isolates from subpopulation 20 (RF: 2.06, 99% CrI: 1.122
388 – 3.778) (Fig 6C). There was an association between subpopulation 18 isolates and older patients
389 (median age 56.1 years) when compared to subpopulation 20 isolates (median age 36.4 years) (p -
390 value: $5.00E-6$, Mann-Whitney U test) (Fig 6D).

391 We hypothesized that the approximately 300kb difference between the assembly lengths
392 of the *S. Infantis* subpopulations may be due to the presence of the pESI plasmid previously
393 identified in *S. Infantis*(36). After checking all isolates for the presence of this plasmid, 144 out of
394 145 *S. Infantis* isolates annotated to subpopulation 18 and 0 out of 243 isolates from

395 subpopulation 20 were putatively positive for pESI plasmids. Only one isolate, a *S. Muenchen*,
396 was putatively positive for the pESI plasmid outside of the *S. Infantis* 18 subpopulation.

397 Two subpopulations represented approximately 85% of the total *S. Typhimurium*
398 population in the analysis set, which we analyzed further (Fig 7A). Similar to the *S. Infantis*
399 subpopulations, the two subpopulations yielded significantly different genome assembly lengths
400 (subpopulation 2: 4.90 Mb, subpopulation 16: 4.85 Mb, p-value < 2.2E-16, Mann-Whitney U
401 test) (Fig 7B). However, the assembly difference of approximately 5kb between the *S.*
402 *Typhimurium* subpopulations is far less dramatic than the approximately 300kb difference
403 observed between *S. Infantis* subpopulations. 668 of the 937 *S. Typhimurium* isolates in
404 subpopulations 2 (n = 359) and 16 (n = 309) have clinical case presentation data. Subpopulation
405 2 isolates presented as double the extraintestinal infections than subpopulation 16 isolates (RF
406 2.11, 99% CrI: 1.109 – 4.016) (Fig 7C). In contrast with the *S. Infantis* subpopulations, the age of
407 patients was not significantly different between the two subpopulations (p-value 0.97, Mann-
408 Whitney U test) (Fig 7D).

409

410 **Fig 6. Description of two *S. Infantis* virulence subpopulations.** (A) Dendrogram highlighting
411 the locations of the two *S. Infantis* virulence subpopulations within the greater population of
412 12,337 *S. enterica* isolates. (B) Histograms of the assembly lengths for the respective
413 subpopulations. (C) Proportion of extraintestinal infections among illnesses caused by the two
414 subpopulations (FoodNet data). (D) Boxplots of the distribution of patient age in infections
415 caused by the two subpopulations (FoodNet data).

416 **Fig 7. Description of two *S. Typhimurium* virulence subpopulations.** (A) Dendrogram
417 highlighting the locations of the two *S. Typhimurium* virulence subpopulations within the greater

418 population of 12,337 *S. enterica* isolates. (B) Histograms of the assembly lengths for the
 419 respective subpopulations. (C) Proportion of extraintestinal infections among illnesses caused by
 420 the two subpopulations (FoodNet data). (D) Boxplots of the distribution of patient age in
 421 infections caused by the two subpopulations (FoodNet data).

422

423 **Table 2. Within serovar virulence subpopulations.**

Serovar	Total Serovar Count	Subpopulation A			Subpopulation B		
		Subpopulation ID	Genome Count	Proportion of Total Serovar Population	Subpopulation ID	Genome Count	Proportion of Total Serovar Population
I 1,4,[5],12:b:-	83	30	30	0.36	29	52	0.63
I 1,4,[5],12:i:-	530	9	294	0.55	2	163	0.31
Infantis	388	20	243	0.63	18	145	0.37
Kentucky	169	32	89	0.53	21	80	0.47
Montevideo	1341	14	503	0.38	11	838	0.62
Muenchen	392	12	133	0.34	6	195	0.50
Newport	1751	6	427	0.24	1	926	0.53
Oranienburg	137	14	32	0.23	11	103	0.75
Reading	60	37	28	0.47	29	28	0.47
Saintpaul	202	29	53	0.26	27	134	0.66
Typhimurium	1106	16	477	0.43	2	460	0.42

424

425 Table 2 legend: Serovars found to contain at least two subpopulations, each representing greater
 426 than 0.20 of the total serovar population. Subpopulations were identified by increasing the
 427 number of clusters to match the number of serovars ($k = 37$). Provided is the subpopulation ID's,
 428 the number of genomes resident within each subcluster, and the proportion of the total
 429 population the subcluster represents. Note, that the subpopulations may not represent the total
 430 combined population of the serovar in the analysis set.

431

432 Discussion

433 The pathogenesis of *S. enterica* is only partially understood, and how different serovars
434 generate distinct disease pathologies is also not well-defined. To better understand how serovars
435 group together based on virulence factor gene carriage, we describe a novel methodology that
436 allowed for rapid identification of serovars of public health concern. Compared to methods used
437 in previous studies [23,24], this scalable genomic approach allowed us to generate a measure of
438 relatedness for a large number of *S. enterica* isolates in a computationally efficient manner and
439 group them using established hierarchical clustering methods [41]. While we considered other
440 clustering methods such as logistic principal component analysis and k-means clustering, we
441 chose the unsupervised random forest approach because it is more robust to outliers, non-
442 parametric, and aggregates results from many models rather than basing inference on a single,
443 “best” model. Our method cannot be read as a traditional phylogeny of evolutionary process but
444 rather as a snapshot of the current virulence potential of more than 12,000 isolates retrieved from
445 humans, bovine animals, and beef products. We did not employ a traditional phylogeny because
446 we were chiefly interested in the current state of potential virulence that consumers are exposed
447 to through beef products, rather than in the evolutionary development of such virulence. As such,
448 we wanted to capture the largest possible number of isolates in a computationally-efficient
449 manner.

450 We contend that this method is pertinent to virulence loci found with SPI as the regions
451 are subject to horizontal gene transfer [45,46]. Common methods to differentiate serovars
452 typically rely on the alignment of core genes or single nucleotide polymorphisms (SNP)
453 identified against reference assemblies [6,24,47,48]. These methods must rely on *post hoc*
454 analysis to determine if two evolutionary similar strains have acquired virulence factor genes

455 which may correspond to differences in case presentation as witnessed in *S. Infantis* and *S.*
456 *Typhimurium*. Demarcating isolates by the presence/absence of virulence factors identified a
457 cluster of higher virulence serovars that accounts for a large proportion of sporadic cases, beef-
458 associated, and total foodborne outbreaks compared to the lower virulence cluster. The higher
459 occurrence of beef-associated outbreaks occurs despite a much lower frequency of isolation from
460 regulatory beef samples relative to serovars from the other clusters combined. Our method, in
461 combination with quantitative risk assessment techniques could be used to account for the
462 relative exposure to serovars (e.g., via different food consumption) and the resultant probability
463 of disease.

464 FoodNet isolates are the basis of the salmonellosis incidence calculation, and the dataset
465 does not provide source attribution. Therefore, serovars from the sporadic human clinical data
466 likely result from multiple exposure sources (poultry, beef, vegetables, etc.). Despite this
467 potentially diverse source of *S. enterica* isolates, most serovars resided in one of the five major
468 clusters (including beef and bovine isolates) suggesting that basal virulence factor gene carriage
469 is conserved within serovars across sources.

470 Interestingly, our method identified a group of virulence factors involved in attachment to
471 host cells or outer membrane structure as the most differential genes between virulence groups
472 (Fig 5). Further, many of these adhesion operons originated from Enterobacteriaceae other than
473 *Salmonella*. Use of only putative *Salmonella* virulence factors from PATRIC [32] and the
474 Virulence Factor Database [31] would not have annotated the open reading frames highlighting
475 the need for expanding the putative virulence factors of *S. enterica* outside of the genera to
476 members of Enterobacteriaceae. In contrast, we found two differentiating operons, *lpf* (long polar
477 fimbriae) and *rfb*, from different virulence categories in significantly higher proportions in

478 cluster 1. The higher proportion of *lpf* genes, which fall into the adhesion virulence category, is
479 notable as the operon has been associated with *S. enterica* binding the Peyer's patches, namely
480 the M-cells found within the lymphoid organs [49,50]. This may potentiate the higher infectivity
481 of cluster 1 serovars as recent work with the Type Three Secretion systems (T3SS) of *S. enterica*
482 (involved with the introduction of effector proteins to the cytoplasm of host cells) suggest that
483 the structure does not penetrate the cytoplasmic membrane like a syringe, but requires tension
484 and adopts a “tent-pole” like structure [51]. If tension is required for the function of the T3SS,
485 enhanced binding to M-cells mediated by the *lpf* operon may be one reason cluster 1 has a higher
486 incidence rate of domestically acquired sporadic cases. The *lpf* operon has not only been
487 implicated in *S. enterica* infections; strains of *Escherichia coli* O157:H7 with mutations in the *lpf*
488 operon show decreased attachment and colonization in both in vitro [52] and in vivo [53] models
489 again show the interplay of virulence mechanisms between Enterobacteriaceae. The roles of *rfb*
490 genes, which were primarily classified in the survival/host persistence virulence category, are not
491 as well investigated as the *lpf* operon but appear to be involved in the biosynthesis of the O-
492 antigen and lipopolysaccharide structuring. A recent report suggests that the full complement of
493 *rfb* genes leads to higher virulence in experimentally infected chickens [54].

494 Examining virulence gene catalogues not only identified large, serovar level clusters but
495 also, by altering the cluster number (k value), virulence subpopulations within serovars. With the
496 current method, it cannot be ascertained whether the virulence subpopulations represent
497 polyphyletic clades within serovars as it cannot be interpreted as a phylogeny. However, by
498 applying a top-down approach, the presence of increased virulence capacity can be readily
499 identified. The two subpopulations of *S. Infantis* present over a two-fold difference in probability
500 of extraintestinal infections. *S. Infantis* has been rapidly increasing in incidence in Israel and

501 previous studies suggest that the addition of a virulence megaplasmid pESI could be responsible
502 [55]. The mean difference between the two subpopulations was approximately 300kb, similar in
503 length to the pESI plasmid (280 kb), and querying *S. Infantis* isolates against a database of
504 marker genes revealed that isolates in the subpopulation with longer assemblies are putatively
505 positive for pESI presence. In addition, the pESI plasmid carries genes necessary for the
506 synthesis of yersiniabactin [55], a siderophore dependent iron uptake system commonly observed
507 in *Yersinia pestis*. The eight genes comprising the *ybt* operon are resident in every strain of the
508 higher invasive cluster of *S. Infantis*, and only two out of 243 isolates from the lower invasive
509 cluster contain the operon. Iron is an essential nutrient for *S. enterica* replication during systemic
510 infections [56]. A previous study suggests that co-infections of Malaria and *S. enterica* leads to
511 more systemic infections as excess iron is released upon the lysis of red blood cells, liberating
512 the metal for use by *S. enterica* [57]. Increased iron availability, due to the addition of
513 yersiniabactin, may be one factor for the almost double rate of extraintestinal infections of *S.*
514 *Infantis* cluster 18 compared to *S. Infantis* infections without this plasmid.

515 The methods employed here cannot identify virulence changes due to sequence variations
516 within virulence loci. Variants of the *macA* and *macB* genes in African strains of *S.*
517 Typhimurium sequence-type 313 may have higher invasiveness in human patients and increased
518 survival against challenge with antimicrobial peptides [58]. Others have identified virulence gene
519 alleles that may correspond to pathogenicity differences [59]. The method employed identifies
520 virulence genes against a non-redundant database using BLASTP, so alleles with variation less
521 than 10% sequence identity will be collapsed into the same gene annotation. Furthermore, we did
522 not consider pseudogene formation of virulence genes. Previous work suggests that pseudogenes
523 in *S. enterica* genomes do not follow neutral evolution (random genetic drift, as in many

524 Eukaryotes) but are readily lost from the chromosome [60]. However, pseudogene formation of
525 the *sseI/srfH* secreted effector protein leads to hyperdissemination of ST313 *S. Typhimurium* in
526 experimentally infected mice [61]. The role of pseudogene formation and the pathogenesis needs
527 more study, and the addition of pseudogene information could further improve virulence
528 classifications. Additionally, we chose to focus our analysis on human, bovine, and beef isolates
529 from the US. It is probable given the diversity of *S. enterica* that all virulence patterns and
530 serovar subpopulations are not represented in this work.

531 *S. enterica* is a diverse pathogen. Yet, most risk assessments and food safety regulations
532 informed by these assessments only separate Typhoidal Salmonellosis and non-Typhoidal
533 Salmonellosis, treating serovars as a homogenous unit [62-64]. Our results suggest that strains
534 with the highest incidence of domestically acquired sporadic cases and outbreaks of human
535 infections share a common virulence repertoire. Control and surveillance programs devoting
536 more resources to clinically relevant serovars might result in increased public health gains, but
537 such interventions must be evaluated using quantitative risk assessment methods.

538 Furthermore, serovar virulence cannot be considered homogenous in all cases as
539 observed with *S. Infantis* and *S. Typhimurium*. Although attributing virulence to specific genes
540 was beyond the scope of this study, our analysis could inform further research to identify
541 *Salmonella* genes associated with severe illness.

542 **Acknowledgements**

543 G.F, J.P., D.T., S.C., and F.Z. are employed by EpiX Analytics. R.P. is a Senior Scientific
544 Advisor for EpiX Analytics. This work utilized the RMACC Summit supercomputer, which is
545 supported by the National Science Foundation (awards ACI-1532235 and ACI-1532236), the
546 University of Colorado Boulder, and Colorado State University. The Summit supercomputer is a

547 joint effort of the University of Colorado Boulder and Colorado State University. FoodNet Data:
548 The findings and conclusions in this report are those of the author(s) and do not necessarily
549 represent the official position of the Centers for Disease Control and Prevention.

550 **References**

- 551 1. Srikantiah P, Lay JC, Hand S, Crump JA, Campbell J, Van Duyn MS, et al. Salmonella
552 enterica serotype Javiana infections associated with amphibian contact, Mississippi, 2001.
553 *Epidemiol Infect.* 2004 Apr;132(2):273–81.
- 554 2. Lawson B, de Pinna E, Horton RA, Macgregor SK, John SK, Chantrey J, et al.
555 Epidemiological Evidence That Garden Birds Are a Source of Human Salmonellosis in
556 England and Wales. *PLOS ONE.* 2014 Feb 26;9(2):e88968.
- 557 3. Salmonella Subcommittee of the Nomenclature Committee of the International Society for
558 Microbiology. The Genus Salmonella Lignières, 1900. *J Hyg (Lond).* 1934 Oct;34(3):333–
559 50.
- 560 4. Grimont P, Weill FX. Antigenic Formulae of the Salmonella serovars, (9th ed.) Paris: WHO
561 Collaborating Centre for Reference and Research on Salmonella. *Inst Pasteur.* 2007 Jan 1;1–
562 166.
- 563 5. Centers for Disease Control and Prevention (CDC). National Salmonella Surveillance
564 Annual Report, 2016. Atlanta, Georgia: US Department of Health and Human Services:
565 CDC; 2018.
- 566 6. Worley J, Meng J, Allard MW, Brown EW, Timme RE. Salmonella enterica Phylogeny
567 Based on Whole-Genome Sequencing Reveals Two New Clades and Novel Patterns of
568 Horizontally Acquired Genetic Elements. *mBio.* 2018 Nov 27;9(6):e02303-18.
- 569 7. Rivera-Chávez F, Bäumlér AJ. The Pyromaniac Inside You: Salmonella Metabolism in the
570 Host Gut. *Annu Rev Microbiol.* 2015 Oct 15;69(1):31–48.
- 571 8. Thiennimitr P, Winter SE, Winter MG, Xavier MN, Tolstikov V, Huseby DL, et al. Intestinal
572 inflammation allows Salmonella to use ethanolamine to compete with the microbiota. *Proc*
573 *Natl Acad Sci U S A.* 2011/10/03 ed. 2011 Oct 18;108(42):17480–5.
- 574 9. Drumo R, Pesciaroli M, Ruggeri J, Tarantino M, Chirullo B, Pistoia C, et al. Salmonella
575 enterica Serovar Typhimurium Exploits Inflammation to Modify Swine Intestinal
576 Microbiota. *Front Cell Infect Microbiol [Internet].* 2016;5. Available from:
577 <https://www.frontiersin.org/article/10.3389/fcimb.2015.00106>
- 578 10. Marcus SL, Brumell JH, Pfeifer CG, Finlay BB. Salmonella pathogenicity islands: big
579 virulence in small packages. *Microbes Infect.* 2000 Feb;2(2):145–56.

- 580 11. Lorkowski M, Felipe-López A, Danzer CA, Hansmeier N, Hensel M. Salmonella enterica
581 invasion of polarized epithelial cells is a highly cooperative effort. *Infect Immun*. 2014/04/07
582 ed. 2014 Jun;82(6):2657–67.
- 583 12. Steele-Mortimer O, Brumell JH, Knodler LA, Méresse S, Lopez A, Finlay BB. The invasion-
584 associated type III secretion system of Salmonella enterica serovar Typhimurium is
585 necessary for intracellular proliferation and vacuole biogenesis in epithelial cells. *Cell*
586 *Microbiol*. 2002 Jan 1;4(1):43–54.
- 587 13. Kurtz JR, Goggins JA, McLachlan JB. Salmonella infection: Interplay between the bacteria
588 and host immune system. *Immunol Lett*. 2017/07/15 ed. 2017 Oct;190:42–50.
- 589 14. Gal-Mor O, Boyle EC, Grassl GA. Same species, different diseases: how and why typhoidal
590 and non-typhoidal Salmonella enterica serovars differ. *Front Microbiol* [Internet]. 2014;5.
591 Available from: <https://www.frontiersin.org/article/10.3389/fmicb.2014.00391>
- 592 15. Harvey RR, Friedman CR, Crim SM, Judd M, Barrett KA, Tolar B, et al. Epidemiology of
593 Salmonella enterica Serotype Dublin Infections among Humans, United States, 1968–2013.
594 *Emerg Infect Dis*. 2017 Sep;23(9):1493–501.
- 595 16. Ramachandran G, Panda A, Higginson EE, Ateh E, Lipsky MM, Sen S, et al. Virulence of
596 invasive Salmonella Typhimurium ST313 in animal models of infection. *PLoS Negl Trop*
597 *Dis*. 2017 Aug 4;11(8):e0005697.
- 598 17. Jiang L, Wang P, Song X, Zhang H, Ma S, Wang J, et al. Salmonella Typhimurium
599 reprograms macrophage metabolism via T3SS effector SopE2 to promote intracellular
600 replication and virulence. *Nat Commun*. 2021 Feb 9;12(1):879.
- 601 18. Cheng RA, Eade CR, Wiedmann M. Embracing Diversity: Differences in Virulence
602 Mechanisms, Disease Severity, and Host Adaptations Contribute to the Success of
603 Nontyphoidal Salmonella as a Foodborne Pathogen. *Front Microbiol* [Internet]. 2019;10.
604 Available from: <https://www.frontiersin.org/article/10.3389/fmicb.2019.01368>
- 605 19. Faber F, Thiennimitr P, Spiga L, Byndloss MX, Litvak Y, Lawhon S, et al. Respiration of
606 Microbiota-Derived 1,2-propanediol Drives Salmonella Expansion during Colitis. *PLOS*
607 *Pathog*. 2017 Jan 5;13(1):e1006129.
- 608 20. Hannemann S, Galán JE. Salmonella enterica serovar-specific transcriptional reprogramming
609 of infected cells. *PLOS Pathog*. 2017 Jul 24;13(7):e1006532.
- 610 21. Pulford CV, Perez-Sepulveda BM, Canals R, Bevington JA, Bengtsson RJ, Wenner N, et al.
611 Stepwise evolution of Salmonella Typhimurium ST313 causing bloodstream infection in
612 Africa. *Nat Microbiol*. 2021 Mar 1;6(3):327–38.
- 613 22. Ebel ED, Williams MS, Schlosser WD. Estimating the Type II error of detecting changes in
614 foodborne illnesses via public health surveillance. *Microb Risk Anal*. 2017 Dec 1;7:1–7.

- 615 23. Karanth S, Tanui CK, Meng J, Pradhan AK. Exploring the predictive capability of advanced
616 machine learning in identifying severe disease phenotype in *Salmonella enterica*. *Food Res*
617 *Int.* 2022 Jan 1;151:110817.
- 618 24. Chen R, Cheng RA, Wiedmann M, Orsi RH. Development of a Genomics-Based Approach
619 To Identify Putative Hypervirulent Nontyphoidal *Salmonella* Isolates: *Salmonella enterica*
620 Serovar Saintpaul as a Model. *mSphere*. 2022 Feb 23;7(1):e0073021.
- 621 25. National Advisory Committee on Microbiological Criteria for Foods (NACMCF). Response
622 to Questions Posed by the Food Safety and Inspection Service: Enhancing *Salmonella*
623 Control in Poultry Products. 2022. Available from:
624 [https://www.fsis.usda.gov/sites/default/files/media_file/documents/NACMCF_Salmonella-](https://www.fsis.usda.gov/sites/default/files/media_file/documents/NACMCF_Salmonella-Poultry_Response_for_Committee_Review.pdf)
625 [Poultry_Response_for_Committee_Review.pdf](https://www.fsis.usda.gov/sites/default/files/media_file/documents/NACMCF_Salmonella-Poultry_Response_for_Committee_Review.pdf)
- 626 26. Ward C. Vertical Integration Comparison: Beef, Pork, and Poultry. Western Agricultural
627 Economics Association; 1997. Available from:
628 <https://EconPapers.repec.org/RePEc:ags:waeare:35759>
- 629 27. Palma F, Manfreda G, Silva M, Parisi A, Barker DOR, Taboada EN, et al. Genome-wide
630 identification of geographical segregated genetic markers in *Salmonella enterica* serovar
631 Typhimurium variant 4,[5],12:i:-. *Sci Rep*. 2018 Oct 15;8(1):15251.
- 632 28. Fenske GJ, Thachil A, McDonough PL, Glaser A, Scaria J. Geography Shapes the
633 Population Genomics of *Salmonella enterica* Dublin. *Genome Biol Evol*. 2019 Aug
634 1;11(8):2220–31.
- 635 29. CDC. Foodborne Diseases Active Surveillance Network. Atlanta, Georgia: US Department
636 of Health and Human Services: CDC;
- 637 30. Yoshida CE, Kruczkiewicz P, Laing CR, Lingohr EJ, Gannon VPJ, Nash JHE, et al. The
638 *Salmonella* In Silico Typing Resource (SISTR): An Open Web-Accessible Tool for Rapidly
639 Typing and Subtyping Draft *Salmonella* Genome Assemblies. *PLOS ONE*. 2016 Jan
640 22;11(1):e0147101.
- 641 31. Liu B, Zheng D, Jin Q, Chen L, Yang J. VFDB 2019: a comparative pathogenomic platform
642 with an interactive web interface. *Nucleic Acids Res*. 2019 Jan 8;47(D1):D687–92.
- 643 32. Mao C, Abraham D, Wattam AR, Wilson MJC, Shukla M, Yoo HS, et al. Curation,
644 integration and visualization of bacterial virulence factors in PATRIC. *Bioinformatics*. 2015
645 Jan 15;31(2):252–8.
- 646 33. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation
647 sequencing data. *Bioinformatics*. 2012/10/11 ed. 2012 Dec 1;28(23):3150–2.
- 648 34. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. 2014 Jul
649 15;30(14):2068–9.

- 650 35. Hyatt D, Chen GL, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
651 gene recognition and translation initiation site identification. *BMC Bioinformatics*. 2010 Mar
652 8;11(1):119.
- 653 36. McMillan EA, Wasilenko JL, Tagg KA, Chen JC, Simmons M, Gupta SK, et al. Carriage
654 and Gene Content Variability of the pESI-Like Plasmid Associated with *Salmonella* *Infantis*
655 Recently Established in United States Poultry Production. *Genes*. 2020 Dec 18;11(12):1516.
- 656 37. Franco A, Leekitcharoenphon P, Feltrin F, Alba P, Cordaro G, Iurescia M, et al. Emergence
657 of a Clonal Lineage of Multidrug-Resistant ESBL-Producing *Salmonella* *Infantis*
658 Transmitted from Broilers and Broiler Meat to Humans in Italy between 2011 and 2014.
659 *PLoS ONE*. 2015 Dec 30;10(12):e0144802.
- 660 38. Seemann T. ABRicate. 2022. [cited 2 December 2022]. Database: github [Internet].
661 Available from: <https://github.com/tseemann/abricate>
- 662 39. Liaw A, Wiener M. Classification and Regression by randomForest. *R News*. 2002;2(3):18–
663 22.
- 664 40. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna,
665 Austria: R Foundation for Statistical Computing; 2021. Available from: [https://www.R-](https://www.R-project.org/)
666 [project.org/](https://www.R-project.org/)
- 667 41. Ward JH. Hierarchical Grouping to Optimize an Objective Function. *J Am Stat Assoc*.
668 1963;58(301):236–44.
- 669 42. Hennig C. Cluster-wise assessment of cluster stability. *Comput Stat Data Anal*. 2007 Sep
670 15;52(1):258–71.
- 671 43. Hennig C. fpc: Flexible Procedures for Clustering. 2020. Available from: [https://CRAN.R-](https://CRAN.R-project.org/package=fpc)
672 [project.org/package=fpc](https://CRAN.R-project.org/package=fpc)
- 673 44. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian Data Analysis*.
674 3rd ed. Boca Raton: CRC Press; 2013.
- 675 45. Lerminiaux NA, MacKenzie KD, Cameron ADS. *Salmonella* Pathogenicity Island 1 (SPI-1):
676 The Evolution and Stabilization of a Core Genomic Type Three Secretion System.
677 *Microorganisms*. 2020 Apr 16;8(4):576.
- 678 46. Brown EW, Bell R, Zhang G, Timme R, Zheng J, Hammack TS, et al. *Salmonella* Genomics
679 in Public Health and Food Safety. *EcoSal Plus*. 2021 Dec 15;9(2):eESP00082020.
- 680 47. Pornsukarom S, van Vliet AHM, Thakur S. Whole genome sequencing analysis of multiple
681 *Salmonella* serovars provides insights into phylogenetic relatedness, antimicrobial resistance,
682 and virulence markers across humans, food animals and agriculture environmental sources.
683 *BMC Genomics*. 2018 Nov 6;19(1):801.

- 684 48. Petrovska L, Mather AE, AbuOun M, Branchu P, Harris SR, Connor T, et al. Microevolution
685 of Monophasic Salmonella Typhimurium during Epidemic, United Kingdom, 2005-2010.
686 *Emerg Infect Dis.* 2016 Apr;22(4):617–24.
- 687 49. Bäumlér AJ, Tsohis RM, Heffron F. The *lpf* fimbrial operon mediates adhesion of Salmonella
688 typhimurium to murine Peyer’s patches. *Proc Natl Acad Sci U S A.* 1996 Jan 9;93(1):279–
689 83.
- 690 50. Gonzales AM, Wilde S, Roland KL. New Insights into the Roles of Long Polar Fimbriae and
691 Stg Fimbriae in Salmonella Interactions with Enterocytes and M Cells. *Infect Immun.* 2017
692 Aug 18;85(9):e00172-17.
- 693 51. Park D, Lara-Tejero M, Waxham MN, Li W, Hu B, Galán JE, et al. Visualization of the type
694 III secretion mediated Salmonella-host cell interface using cryo-electron tomography. *eLife.*
695 2018 Oct 3;7.
- 696 52. Torres AG, Kanack KJ, Tutt CB, Popov V, Kaper JB. Characterization of the second long
697 polar (LP) fimbriae of *Escherichia coli* O157:H7 and distribution of LP fimbriae in other
698 pathogenic *E. coli* strains. *FEMS Microbiol Lett.* 2004 Sep 1;238(2):333–44.
- 699 53. Jordan DM, Cornick N, Torres AG, Dean-Nystrom EA, Kaper JB, Moon HW. Long polar
700 fimbriae contribute to colonization by *Escherichia coli* O157:H7 in vivo. *Infect Immun.* 2004
701 Oct;72(10):6168–71.
- 702 54. Gao R, Huang H, Hamel J, Levesque RC, Goodridge LD, Ogunremi D. Application of a
703 High-Throughput Targeted Sequence AmpliSeq Procedure to Assess the Presence and
704 Variants of Virulence Genes in Salmonella. *Microorganisms.* 2022 Feb 5;10(2).
- 705 55. Aviv G, Tsyba K, Steck N, Salmon-Divon M, Cornelius A, Rahav G, et al. A unique
706 megaplasmid contributes to stress tolerance and pathogenicity of an emergent Salmonella
707 enterica serovar Infantis strain. *Environ Microbiol.* 2014 Apr 1;16(4):977–94.
- 708 56. Nairz M, Ferring-Appel D, Casarrubea D, Sonnweber T, Viatte L, Schroll A, et al. Iron
709 Regulatory Proteins Mediate Host Resistance to Salmonella Infection. *Cell Host Microbe.*
710 2015 Aug 12;18(2):254–61.
- 711 57. van Santen S, de Mast Q, Swinkels DW, van der Ven AJAM. The iron link between malaria
712 and invasive non-typhoid Salmonella infections. *Trends Parasitol.* 2013 May;29(5):220–7.
- 713 58. Honeycutt JD, Wenner N, Li Y, Brewer SM, Massis LM, Brubaker SW, et al. Genetic
714 variation in the *MacAB-TolC* efflux pump influences pathogenesis of invasive Salmonella
715 isolates from Africa. *PLOS Pathog.* 2020 Aug 24;16(8):e1008763.
- 716 59. Rakov AV, Mastriani E, Liu SL, Schifferli DM. Association of Salmonella virulence factor
717 alleles with intestinal and invasive serovars. *BMC Genomics.* 2019 May 28;20(1):429.
- 718 60. Kuo CH, Ochman H. The Extinction Dynamics of Bacterial Pseudogenes. *PLOS Genet.*
719 2010 Aug 5;6(8):e1001050.

- 720 61. Carden SE, Walker GT, Honeycutt J, Lugo K, Pham T, Jacobson A, et al. Pseudogenization
721 of the Secreted Effector Gene *sseI* Confers Rapid Systemic Dissemination of *S.*
722 Typhimurium ST313 within Migratory Dendritic Cells. *Cell Host Microbe*. 2017 Feb
723 8;21(2):182–94.
- 724 62. Food Safety Inspection Service (FSIS). Public Health Effects of Performance Standards for
725 Ground Beef and Beef Manufacturing Trimmings. Washington D.C.: US Department of
726 Agriculture: FSIS; 2019.
- 727 63. Food Safety Inspection Service (FSIS). Public Health Effects of Raw Chicken Parts and
728 Comminuted Chicken and Turkey Performance Standards. Washington D.C.: US
729 Department of Agriculture: FSIS; 2015.
- 730 64. Lambertini E, Ruzante JM, Kowalczyk, BB. The Public Health Impact of Implementing a
731 Concentration-Based Microbiological Criterion for Controlling Salmonella in Ground
732 Turkey. *Risk Analysis*. 2021 Aug; 41(8):1376-95.

733

734 **Supporting information**

735 **S1 Fig. Addition of a sixth virulence cluster.** (A) Dendrogram depicting the hierarchical
736 relationship between 12,337 *S. enterica* genome assemblies based upon virulence factor gene
737 carriage with six virulence clusters superimposed on top. (B) Heatmap of serovar proportion
738 within each of the six respective virulence clusters. Rows are clustered using Ward's method. (C)
739 Characteristics of the six virulence clusters: cluster stability - Jaccard similarity of 10,000 non-
740 parametric bootstraps, Number of Genomes - depicting the number of *S. enterica* genomes
741 constituent in each cluster, and number of serovars (within cluster serovar proportion > 0.5) in
742 each cluster.

743

744 **S1 Table. Metadata for the analysis set of genomes and SISTR serovar prediction.** Metadata
745 for the contig assemblies used in the analysis including results of the in silico serovar prediction
746 for the analysis set genomes from the SISTR software.

747

748 **S2 Table. Serovar virulence cluster designations.** Virulence cluster designations ($k = 5$) for
749 the 37 serovars in the analysis set.

750

751 **S3 Table. Full list of putative virulence loci considered in the random forest model.** Gene
752 name, locus tag, database source, Genus origin, gene product and classification for the 182
753 putative virulence factor loci used in the random forest model.

754

755 **S4 Table. Epidemiological indicators computed for each virulence cluster.** Estimates of:
756 incidence of domestically acquired sporadic cases per 100k people per year, hospitalization
757 proportion given infection, proportion positive in FSIS testing of US beef products (MT43,
758 MT60, MT64), extraintestinal proportion given infection, and mortality proportion given
759 infection.

760

761 **S5 Table. Virulence factors present in all isolates.** List of virulence factor genes ($n=68$)
762 present in all isolates.

763

764 **S6 Table. Isolate virulence subpopulation cluster designations.** Subpopulation cluster
765 designations ($k = 37$) for the 12,337 contig assemblies in the analysis set.

766

767

768

Bovine and Beef *S. enterica* Isolates

Human *S. enterica* Isolates

Isolates from beef (FSIS HACCP sampling and FDA NARMS), and bovine isolates (FSIS NARMS)

n = 38,795

Ineligible: n = 32,285
Isolation sources not from bovine animals nor beef.

Bovine or beef associated isolates

n = 6,510

Human associated isolates from FoodNet & NORS surveillance

n = 15,054

Human associated isolates from US sporadic cases and beef-associated outbreaks.

n = 7,793

Ineligible: n = 7,261

1. FoodNet:
 - A. International travel (n = 974)
 - B. Immigration related (n = 11)
 - C. Outbreak origin (n = 905)
2. NORS: Attribution to source other than beef (n = 4,847)

Analysis Set

Combined set of isolates

n = 14,303

Ineligible: n = 494

1. No assembly (276)
2. Assembler not SKESA v. 2.2 (n = 149)
3. Contig n > 300 (n = 69)
4. Contig n50 < 2.5E4 bp (n = 0)

Isolates passing initial assembly QC

n = 13,809

Ineligible: n = 162

Assembly fails QC step in serovar prediction.

Isolate assemblies with predicted serovar

n = 13,647

Ineligible: n = 1,310

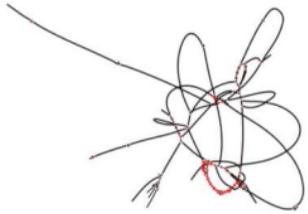
Isolate serovar contains less than 50 isolates

Final Isolate Assembly Set

Bovine or Beef Associated : n = 5,586
Human Associated: n = 6,751

Total: n = 12,337

Acquire and Quality Control Genome Assemblies

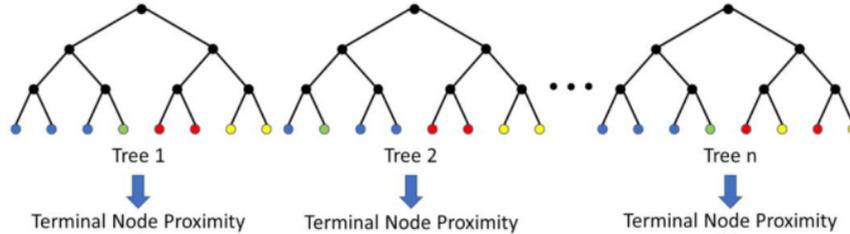


Identify Virulence Genes

Isolate A 111111000100111200101111
 Isolate B 111100011101000211110001
 Isolate C 011100011120001111111101
 Isolate D 111111000100111200100000
 Isolate E 111010000101111200101101
 Isolate F 111111000100111200100111
 Isolate G 111100011101000211110001
 Isolate H 111100011101000211111112

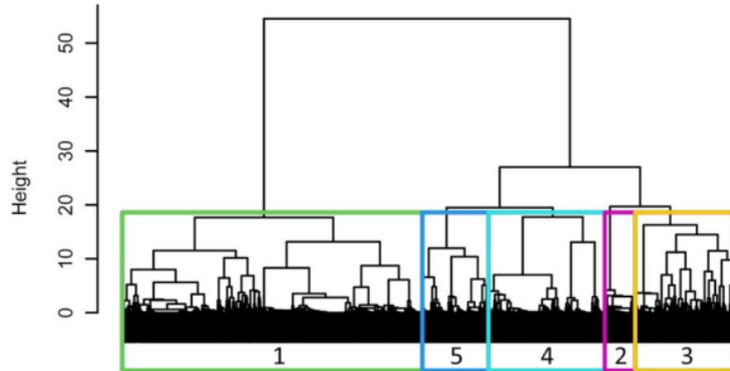


*Compute Isolate Relatedness
(Unsupervised Random Forest)*

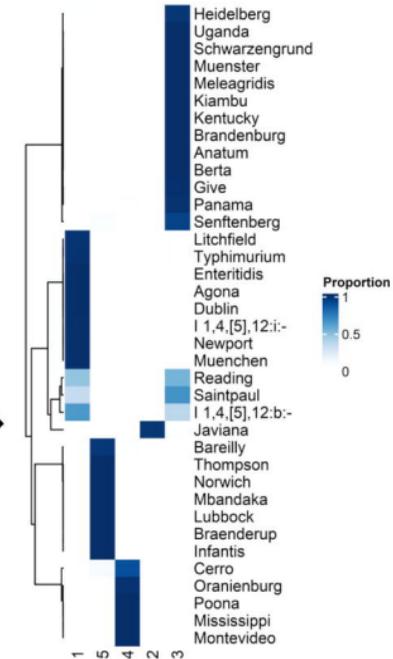


Isolate Cluster Generation

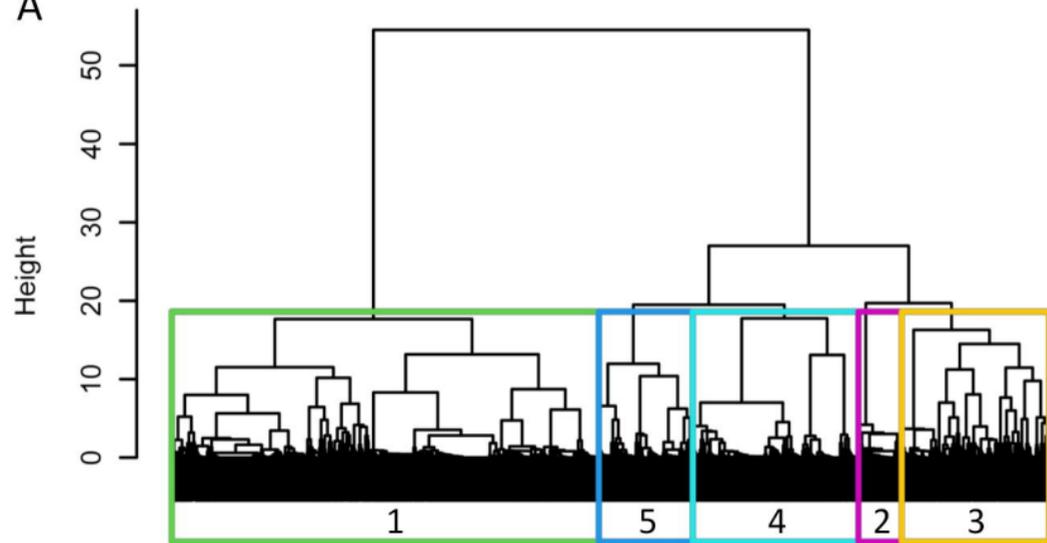
(Hierarchical clustering and non-parametric bootstrapping)



Serovar Virulence Clusters



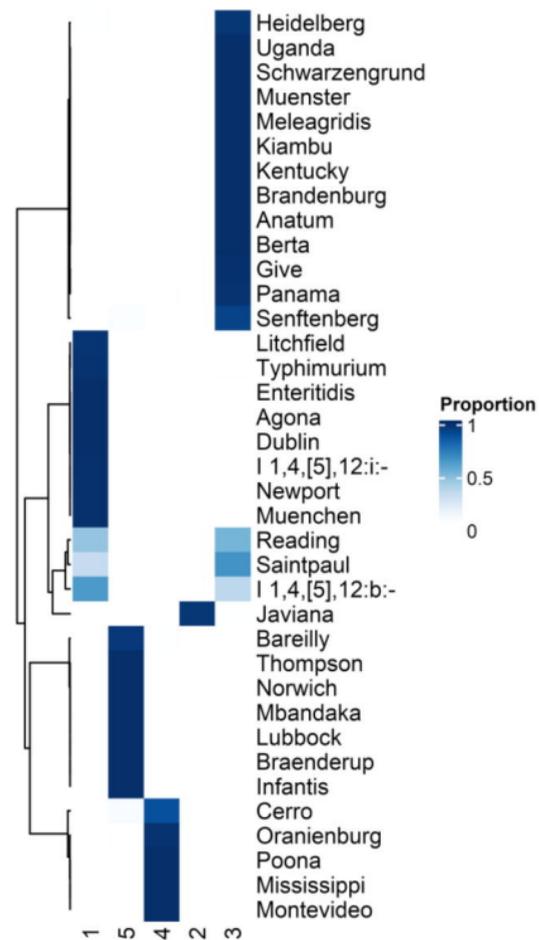
A



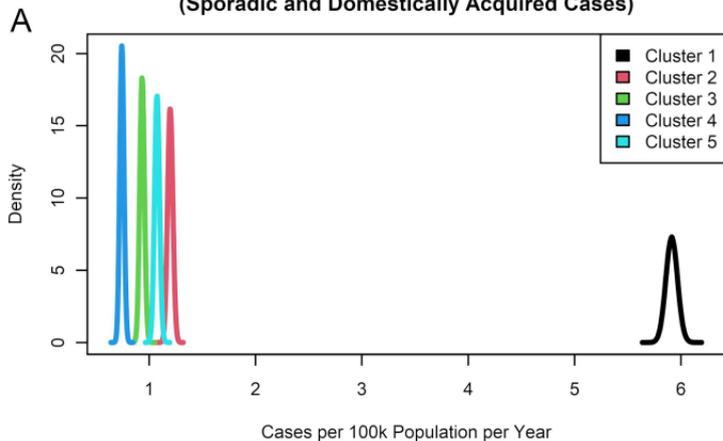
C

Cluster	Cluster Stability (Jaccard Similarity)	Number of Genomes	Serovar Count (prop > 0.50)
1	0.960	6006	9
2	0.887	608	1
3	0.838	2073	15
4	0.995	2329	5
5	0.987	1321	7

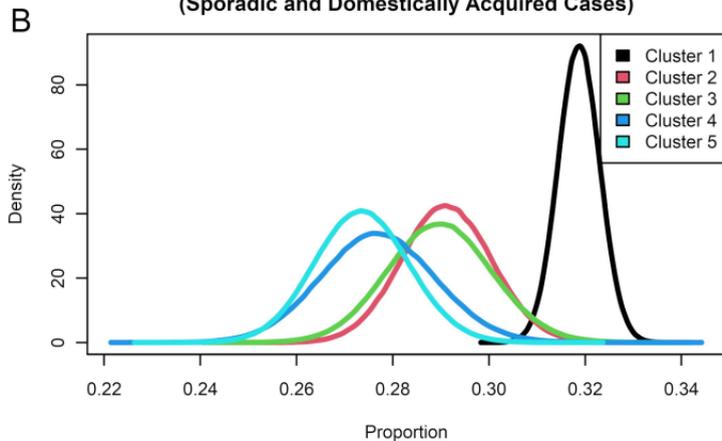
B



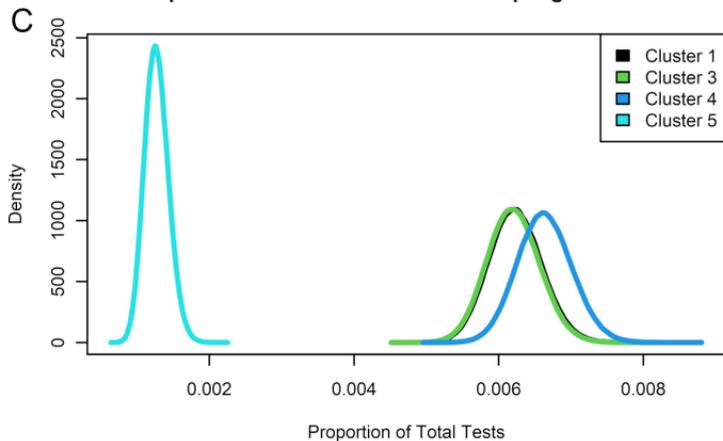
Average Annual Incidence Rate: 2016 - 2019
(Sporadic and Domestically Acquired Cases)



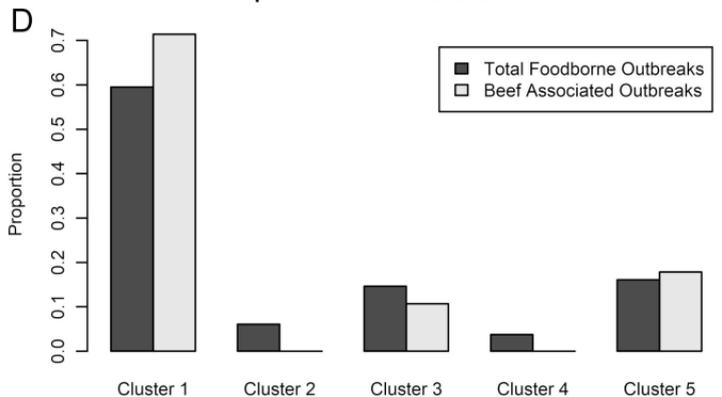
Hospitalization Proportion: 2016 - 2019
(Sporadic and Domestically Acquired Cases)



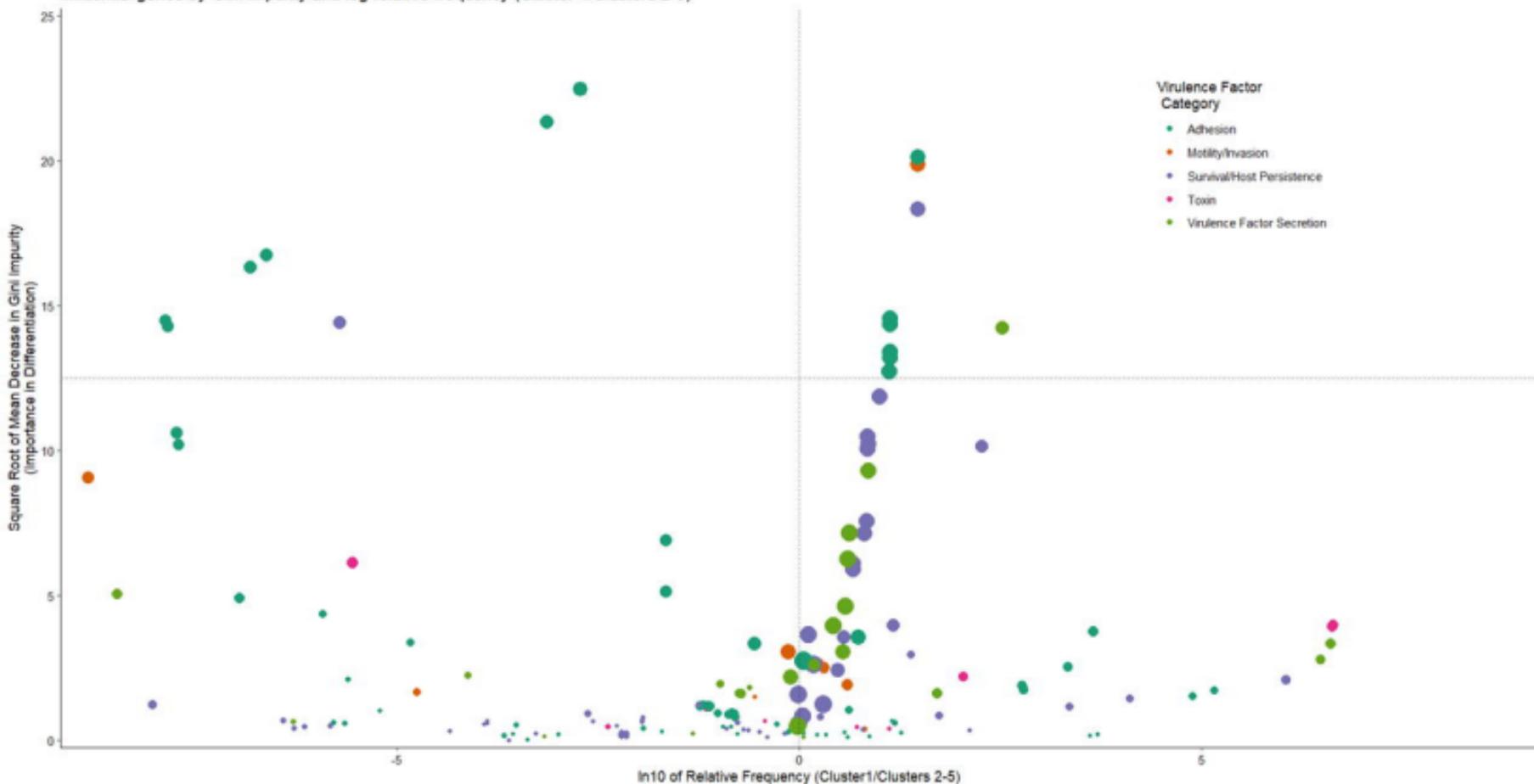
Proportion Positive in FSIS Beef Sampling: 2016-2019



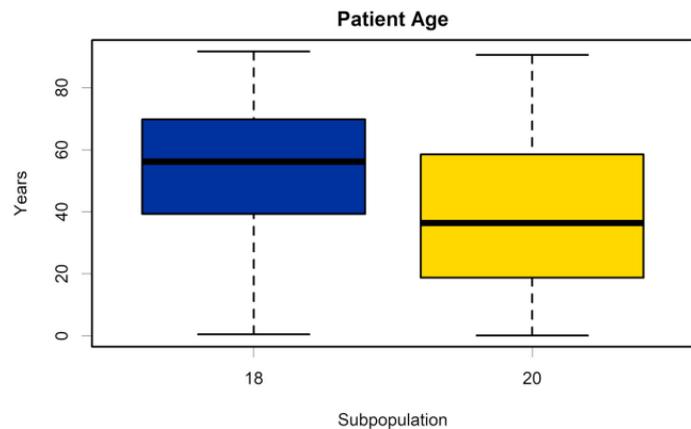
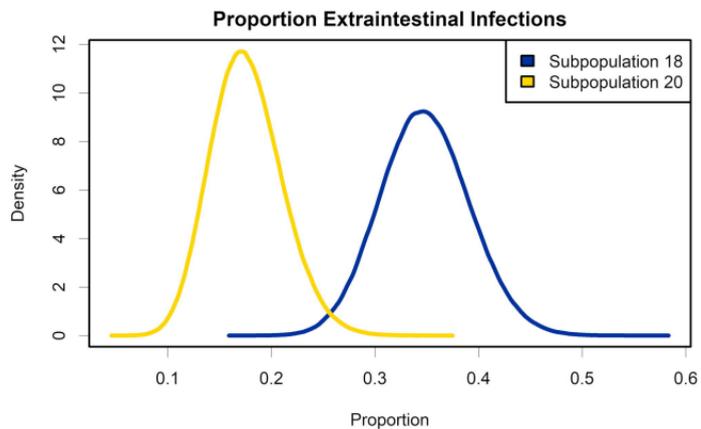
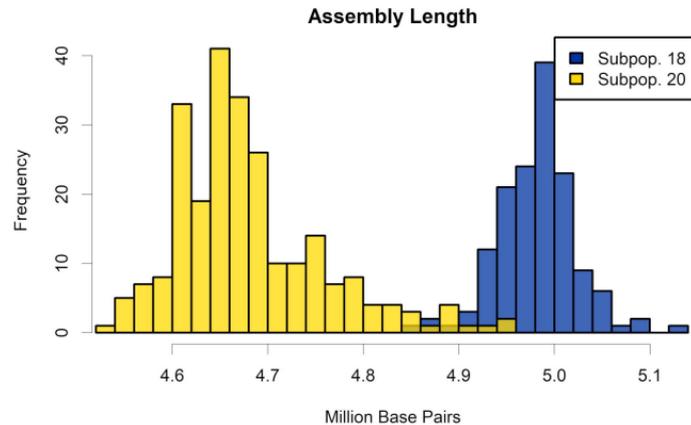
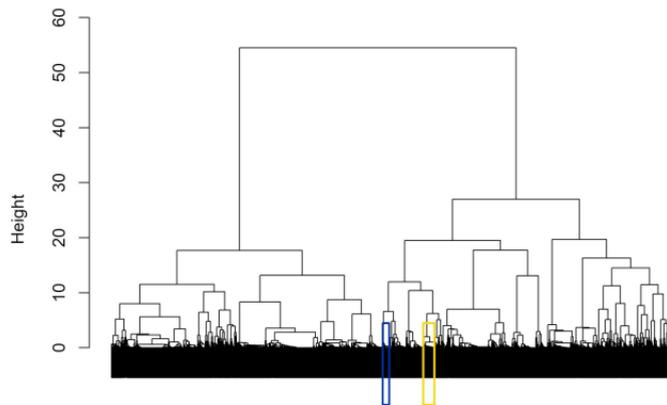
Proportion of Outbreaks: 2016 - 2019



Influential genes by Gini impurity and log relative frequency (Cluster 1/Clusters 2-5)



S. Infantis



S. Typhimurium

