

1

Is it the same strain?

2 **Defining genomic epidemiology thresholds tailored to individual outbreaks**

3

4 **Authors**

5 Audrey Duval, PhD^{1,2,3}, Lulla Opatowski, PhD^{1,2} and Sylvain Brisse, PhD³

6 **Affiliations**

7 ¹Epidemiology and modelling of bacterial escape to antimicrobials, Institut Pasteur, Paris,
8 France.

9 ²Anti-infective Evasion and Pharmacoepidemiology Team, CESP, Université Paris-Saclay,
10 UVSQ, INSERM U1018, Montigny-le-Bretonneux, France.

11 ³ Institut Pasteur, Biodiversity and Epidemiology of Bacterial Pathogens, Paris, France

12

13 **Correspondence**

14 **Sylvain Brisse:** Institut Pasteur, Biodiversity and Epidemiology of Bacterial Pathogens,

15 28 rue du Docteur Roux, F-75724 Paris, France. Phone: +33 1 45 68 83 34 ; E-mail:

16 sylvain.brisse@pasteur.fr

17

18 **Keywords:** genomic epidemiology; outbreak genetic threshold; point-source outbreaks; strain

19 definition

20

Abstract

21 **Background**

22 Epidemiological surveillance relies on microbial strain typing, which defines genomic
23 relatedness among isolates to identify case clusters and their potential sources. No consensus
24 exists on the choice of thresholds of genomic relatedness to define clusters. While *a priori*
25 defined thresholds are often applied, outbreak-specific features such as pathogen mutation
26 rate and duration of source contamination should be considered.

27 **Methods**

28 We developed a forward model of bacterial evolution to simulate mutation within a
29 population diversifying at a specific mutation rate, with specific outbreak duration and sample
30 isolation dates. Based on the resulting expected distribution of genetic distances we define a
31 threshold beyond which isolates are considered as not part of the outbreak. We additionally
32 embedded the model into a Markov Chain Monte Carlo inference framework to estimate,
33 from data including sampling dates or isolates genetic variation, the most credible mutation
34 rate or time since source contamination.

35 **Findings**

36 A simulation study validated the model over realistic durations and mutation rates. When
37 applied to 16 published datasets describing foodborne outbreaks, our framework consistently
38 identified outliers. Appropriate thresholds for grouping cases were obtained for 14 outbreaks.
39 For the remaining two outbreaks, re-estimation of the duration of outbreak lead to updated
40 threshold values and was more likely, given our model, to result in the observed genetic
41 distances.

42 **Interpretation**

43 We propose an evolutionary approach to the ‘single strain’ conundrum by defining the genetic
44 threshold based on individual outbreak properties. The framework provides an informed
45 estimation of the likelihood of a cluster given the samples epidemiological and
46 microbiological context. This forward model, applicable to foodborne or environmental-
47 source single point case clusters or outbreaks, will be useful for epidemiological surveillance
48 and to guide control measures.

49 **Funding**

50 This work was supported financially by the MedVetKlebs project, a component of European
51 Joint Programme One Health EJP, which has received funding from the European Union's
52 Horizon 2020 research and innovation programme under Grant Agreement No 773830. The
53 funders had no role in study design, data collection and analysis, decision to publish, or
54 preparation of the manuscript.

55

56 **Research in context**

57 **Evidence before this study**

58 We searched PubMed for studies published between database inception and April 3, 2021,
59 with the term (threshold OR cut-off OR genetic relatedness) AND (outbreak) AND (cgMLST
60 OR wgMLST OR SNPs) AND (microbial OR bacteria OR bacterial OR pathogen). We found
61 222 related articles. Most studies define a fixed SNP threshold that relate outbreak strains
62 based on previous observations. One original study identifies outbreak clusters based on
63 transmission events. However, it relies on strong assumptions about molecular clock and
64 transmission processes.

65 **Added value of this study**

66 Our study describes a new method based on a forward Wright-Fisher model to find the most
67 credible genetic distance threshold. This method is fast and simple to use with only few
68 assumptions, informed by outbreak duration and pathogen mutation rate. By using SNP or
69 cgMLST pairwise distances and sample collection dates of the outbreak of interest, the
70 algorithm provides context-based guidance to separate outbreak strains from outliers.

71 **Implications of all the available evidence**

72 The fast and easy method developed here enables to move away from *a priori* defined
73 thresholds. Defining clusters more accurately based on the specific features of outbreaks, and
74 the ability to estimate outbreak duration, will provide the needed precision for
75 epidemiological surveillance and should contribute to leverage molecular epidemiology data
76 more efficiently for the purpose of uncovering contamination sources.

77

78 **Data Availability Statement**

79 All data and code used for this manuscript is available online at <https://gitlab.pasteur.fr/BEBP>.

80 **Introduction**

81 Outbreaks of infections caused by the exposure to a unique source are the particular focus of
82 surveillance and infection control strategies. The rapid identification of the source can lead to
83 immediate public health benefits and is therefore critical. In the simplest cases, a single strain
84 of infectious agent contaminates the source and subsequently causes infections (referred to as
85 a ‘clonal outbreak’). This is often the case for contaminated food, water or environmental
86 sources that are under strong regulatory measures and typically uncontaminated. Surveillance
87 systems were therefore put in place, *e.g.*, for food-borne pathogens such as *Salmonella* or
88 *Listeria monocytogenes*, based on a collect-genotype-compare strategy [1,2]. This strategy,
89 dubbed ‘reverse epidemiology’ [3], forms the basis of surveillance systems for foodborne
90 pathogens, such as PulseNet [1]. Molecular surveillance (‘genetic fingerprinting’) enables the
91 detection of nearly identical infectious agent isolates and may trigger epidemiological
92 investigations. These include the search for case-associated risk factors as well as
93 microbiological analyses of suspected sources, and may lead to infection control measures
94 that can prevent further infections.

95 Distinguishing case cluster isolates from sporadic ones has been the ‘Holy Grail’ of molecular
96 epidemiological surveillance. However, the identification of single-strain clusters of
97 infections is confounded by a background of sporadic cases caused by exposure to unrelated
98 sources. Defining ‘a single strain’ typically involves the use of a threshold of genetic distance,
99 which discriminates between isolates that are related or not to the event. The literature is ripe
100 with attempts to define such thresholds [4]. In the whole-genome sequencing (WGS) era,
101 thresholds were refined compared to pre-genomic methods such as PFGE [5–10]. Usually,
102 threshold definition is based on the variability observed within previously well-characterised
103 outbreaks, an approach rooted in the epidemiological concordance principle [11]. However,
104 interpretation of molecular data for strain definition is far from being consensual [5,12,13].

105 From an evolutionary biology point of view, infectious agents that are present as
106 contaminants of an initially sterile source can be considered as subpopulations of individuals
107 that have evolved from a single common ancestor (the original strain) since some time (the
108 duration of contamination). Major factors expected to influence the genetic distances among
109 sampled individuals (isolates) include: i) the duration of strain persistence in the contaminated
110 source prior to infections; ii) the evolutionary rate of the pathogen genomic markers; iii) the
111 sampling dates. On the other hand, the genetic distance to the closest observed isolate

112 unrelated by source will be determined by which genomes were sampled outside the
113 contamination event. All these parameters considered, the quest for a unique threshold
114 applicable to all outbreaks is deemed to fail. Instead, using outbreak-specific thresholds
115 defined based on their context-informed expected diversity is likely to represent a more
116 successful strategy. Attempts to ground threshold definition in evolutionary biology are recent
117 and used the coalescent model [6], transmission models [14] and Bayesian MRCA models
118 [15,16].

119 The aim of this work was the development of a novel model to define the most credible
120 genetic distance cut-offs for single strain outbreaks from a contaminated source, by simulating
121 the accumulation of mutations using specific outbreak parameters.

122 **Methods**

123 **Evolutionary model and definition of the outbreak genetic distance threshold**

124 We define an outbreak (or cluster of cases) as a group of infection cases caused by a single
125 strain ('monoclonal'), excluding co-occurring cases caused by genetically unrelated strains
126 (*i.e.*, from other sources). In the case where two or more genetically unrelated strains co-
127 contaminate the source of the outbreak, they should be analysed separately with this
128 framework.

129 Our evolutionary formalization (**Figure 1A**) is based on a Wright-Fisher forward model of
130 haploid infectious agent evolution [17,18] with constant population size. The simulation is
131 initialised with a homogeneous population of an infectious agent characterised by five
132 properties: i) L , the genome length (base pairs, bp) or the average length of genes of
133 multilocus sequence typing [MLST] approaches; ii) g , the number of genes; iii) μ , the number
134 of substitutions per site per year; iv) D , the duration (in days) of the outbreak, defined as the
135 time elapsed between the initial contamination of the source, and the sampling date of the last
136 isolate; and v) S_d , the set of sampling dates of isolates, which is defined either directly from
137 the source sampling dates or from the date of sampling of infections, in which case the
138 incubation time and within-patient evolution is neglected. Substitutions are introduced at each
139 time step in individuals sampled with replacement according to a uniform distribution
140 (Poisson distribution with parameter λ). A distribution of pairwise genetic distances is
141 generated on these sampled individuals, and the genetic threshold value is defined from this
142 distribution. Details of the model are provided in the **Supplementary Appendix**.

143

144 **Analysis of published outbreak datasets**

145 We reviewed available published outbreak datasets from the literature and analysed the 16
146 datasets listed in **Table 1** [6,19–24] using our modelling framework. Inclusion criteria were i)
147 foodborne outbreak; ii) the availability of whole genome sequence data and iii) availability of
148 collection dates of isolates. The 16 outbreaks are described in more details in the
149 **supplementary appendix**. We estimated D based on evidence provided in the original
150 publications on these outbreaks. We also used previously estimated μ and g for the
151 corresponding infectious agent from literature (**Table 1**). We labelled D and μ values taken

152 from the literature as D_{lit} and μ_{lit} , whereas those derived from our Markov Chain Monte Carlo
153 (MCMC) estimation (see below) were labelled as $D_{estimated}$ and $\mu_{estimated}$.

154

155 **Statistical analyses, simulation studies and statistical framework**

156 *Model assessment.* To assess the capacity of the model to adequately tell apart outbreak
157 isolates from non-outbreak isolates, we used synthetic datasets generated with different
158 parameters values. We applied our framework to a series of 171 simulated outbreaks
159 generated with 19 different values of D each combined with 9 values of μ and including
160 simulated sporadic isolates (**Table S1** in the supplementary appendix). For each of them, we
161 assessed the global sensitivity (Se) and specificity (Sp) of the framework. Details are provided
162 in the **Supplementary Appendix**.

163 *Parameters estimation.* Our model was embedded into a Bayesian inference statistical
164 framework to enable estimation of either the duration (D) or the substitution rate (μ) of
165 studied outbreaks (**Figure 1B; Supplementary appendix**). Simulated outbreaks were used to
166 assess the ability of the model to estimate D and μ , and their impact on the genetic threshold
167 estimation. We used the Kolmogorov-Smirnoff test statistic (noted D_{KS}) to compare real
168 distributions with simulated distributions as a goodness of fit indicator. Details on the
169 inference framework are provided in the **Supplementary Appendix**.

170

171 **Role of the funding source**

172 The funding source did not have an involvement in either study design, collection, analysis, or
173 interpretation of the data.

174 Results

175 Analysis of simulated outbreaks: accuracy of outbreak delineation and of parameters 176 estimation

177 To test the ability of the framework in distinguishing between outbreak and non-outbreak
178 cases, we generated synthetic outbreaks from different combinations of D and μ (**Table S1 in**
179 **the supplementary appendix**). **Figure 2** shows the specificity S_p and sensitivity S_e
180 according to μ . S_p was poor with low μ values, especially when R_d (the ratio of evolution
181 duration between outbreak and non-outbreak genomes) was small (**Figure 2A**). In contrast, S_e
182 was always high (more than 99%, **Figure 2B**), irrespective of the parameter's combinations.
183 We observed that the higher R_d and μ were, the lower this 95% S_p D -value threshold was
184 (**Figure 2C**).

185 We next evaluated whether the model and framework could accurately estimate the
186 parameters D and μ from outbreaks data. To do so, we simulated synthetic outbreaks for
187 which the D and μ values were known, and attempted to estimate one or the other. Regarding
188 D estimation, all HPD include the true value, with higher values of D being associated with
189 smaller 95% HPD (**Figure 3A**). Similarly, μ was adequately estimated, with best estimates
190 being closer to the target value for higher μ values (**Figure 3B**).
191 Because higher D and/or μ values lead in average to more SNPs, we indeed expected more
192 precision in HPDs estimates in these cases.

193 We also investigated the impact of sampling density on estimation accuracy. Results suggest
194 that poor sampling densities (*e.g.*, 5%), when associated to low values of D and μ (therefore
195 resulting in a low genetic diversity among samples), resulted in biased estimations of D and μ ,
196 which were generally overestimated (**Figure 4A** and **4B**). However, we show that sampling
197 densities >10% led to unbiased estimations.

198

199 Genetic threshold definition for published outbreak datasets

200 For each of the 16 published outbreaks, we applied our framework to estimate an expected
201 outbreak-specific genetic threshold value (**Figure 5** provides the example of outbreak 11; see
202 Supplementary appendix figures S1 to S16 for all outbreaks). We found that, for 14 out of 16
203 outbreaks, the classification of isolates as being outbreak-related or sporadic is consistent with

204 previously reported results. Four of these outbreaks included outliers (outbreaks 1, 4, 12 and
205 16), which are correctly classified beyond the threshold of exclusion by our model, except for
206 one isolate of outbreak 4 (**Table 1; Fig S4**; note that outbreak 4 comprised three different co-
207 contaminating genetic clusters [20]; here the defined outbreak strain was ST528). Ten other
208 outbreaks (2, 3, 5, 6, 7, 9, 10, 13, 14 and 15) have no sporadic cases, and our framework
209 correctly clusters all suspected isolates as outbreak-related.

210 For two of the 16 outbreaks, our model leads to different conclusions compared with previous
211 results. In outbreak 8 (*L. monocytogenes*, beef), two isolates are classified as outliers by our
212 model, whereas they were initially classified as outbreak-related [24]. In outbreak 11
213 (*L. monocytogenes*, ox tongue), two isolates came from food and two others from humans.
214 Our algorithm separates food samples in one cluster and human samples in another cluster,
215 whereas the isolates were initially grouped based on epidemiological and genetic evidence:
216 here, the threshold inferred by our model was smaller.

217 When evaluating the influence of outliers on the inferred threshold by removing them from
218 the analysis we find that, in all cases, the outliers do not affect the outbreak threshold. For
219 outbreak 1, 4 and 16, this removal does not change the threshold value but improves the fit
220 between the pairwise SNP distance distribution of the data and the simulated one
221 (**Supplementary Table S2**).

222

223 **Estimation of D and μ values from real outbreaks, and impact on outbreak definition**

224 For each of the 16 above outbreaks, we used our framework to estimate outbreak duration D
225 and substitution rate μ (called $D_{estimated}$ and $\mu_{estimated}$) separately, and used these values (instead
226 of D_{lit} and μ_{lit} used above) for the inference of the genetic distance threshold. Results are
227 provided in **Table 1**.

228 For 11 of the 16 outbreaks, the estimated HPD intervals include D_{lit} . For the 5 remaining, we
229 find higher $D_{estimated}$ values compared with previously reported D_{lit} (**Figure S6 and Table 1**).
230 Regarding μ , for nine outbreaks HPD intervals include their corresponding μ_{lit} , whereas for
231 only one outbreak $\mu_{estimated}$ is lower than μ_{lit} (outbreak 2) and the six remaining outbreaks lead
232 to a higher estimated $\mu_{estimated}$ compared with μ_{lit} . It is important to note that the $D_{estimated}$ 95%
233 HPD is also higher than D_{lit} for these same 6 outbreaks (**Table 1**).

234 After reanalysing the outbreaks using our $D_{estimated}$ and $\mu_{estimated}$ values, we observe that the
235 newly obtained thresholds do not affect the attribution of isolates to the outbreak or sporadic
236 categories in most cases, with three exceptions. First, for outbreak 4, using $D_{estimated}$ or $\mu_{estimated}$
237 increases the threshold from 4 to 11 SNPs, leading to add the previously missing isolate but
238 still excluding the outliers. Second, for outbreak 15, a decreased genetic threshold (4 SNPs
239 instead of 5, in both independent estimations analyses for $D_{estimated}$ and $\mu_{estimated}$) leads to the
240 exclusion of one isolate. Third, for outbreak 11, the genetic threshold is increased from 4
241 SNPs to 7 and 10 SNPs (using $D_{estimated}$ and $\mu_{estimated}$ respectively), leading to group all isolates
242 from food and human samples (**Figure 5**). We also observe that in most cases, using the
243 estimated values of D and μ improves the fit of the genetic distance distribution, with two
244 exceptions (**Table S2** in the supplementary appendix).

245 Discussion

246 Molecular surveillance contributes to identify common exposure to a specific source even
247 when dates and places of infections are distant [25–27]. Given the large differences existing
248 among outbreaks, it is being increasingly recognised that no single-species threshold can be
249 applied to distinguish between outbreak and non-outbreak isolates. To our knowledge,
250 Octavia and colleagues (2015) were the first to attempt to model the expected genetic distance
251 among food outbreak isolates. Although the authors incorporated mutation rate and outbreak
252 duration in their model, they did not use the actual sampling dates. Consequently, their
253 proposed thresholds depend on strong assumptions as to the actual duration of the outbreak
254 (referred to as the *ex-vivo/in-vivo* evolution time by these authors). Stimson *et al.* [14]
255 modelled the number of transmissions that separates infection cases, using a probabilistic
256 model that incorporates the transmission process in addition to mutation rate and timing of
257 infections. Because it models between-host transmission, this approach does not apply to
258 point-source food outbreaks. Lastly, Coll *et al.* [28] aimed at defining a SNP threshold above
259 which transmission of *S. aureus* between humans can be ruled out, by incorporating the
260 timing of transmission and within-host diversity. This evolutionary modelling approach
261 provides a robust SNP cut-off applicable to this specific ecological situation.

262 We propose an original evolutionary approach to the ‘single strain’ threshold conundrum by
263 incorporating epidemiological and microbiological specifics of each outbreak. Our model is
264 supported by a high sensitivity (>90%) of isolates classification and by the results of analyses
265 of 16 real-life published datasets from foodborne outbreaks, which led to consistent results in
266 most cases and enabled to refine outbreak analysis in two cases.

267 The simulation study showed that our model performed well at grouping outbreak cases. We
268 also observed that as D and μ increased, the estimated genetic threshold was more accurate:
269 the model specificity increased with genetic diversity. This is akin to higher resolution typing
270 methods being better at discriminating related and non-related cases. We also found an impact
271 of the evolutionary distance between outbreak and sporadic isolates on model specificity,
272 consistent with the known uncertainty in ruling out sporadic cases for genetically
273 homogeneous pathogens. In addition, we found that the sampling density is important, as it
274 influences the number of observed genetic differences: outbreaks with low diversity will
275 require more samples to capture enough pairwise differences for estimation purposes.

276 Our model assumes a constant population, to avoid increasing execution time with growing
277 bacterial populations. Because the population N remains constant over time, this number must
278 be chosen high enough to capture all the diversity through our sampling process. Indeed, we
279 simulated the sampling processes and did not analyse the whole N population. Because λ , the
280 Poisson parameter, is defined as a function of N , a number of 500 or 1000 is usually enough
281 to capture all bacterial diversity, but higher values should be tested further when extreme
282 substitution rates or duration are explored.

283 In most outbreak investigations, the time since source contamination is unknown, and the
284 underestimation of D is a common risk given the possibility of cryptic transmission and
285 unreported cases having occurred prior to actual outbreak detection [29]. Prior knowledge of
286 μ is also subjected to uncertainty: this parameter strongly depends on the species but also on
287 the strain [30] and on other conditions (*e.g.* temperature, cellular stress). We showed that,
288 although the estimates were largely consistent with epidemiological information, estimated D
289 and μ were often larger. As D and μ both affect the expected genetic diversity in the same
290 direction, it is impossible to decide whether it is the rate, or the duration, that was higher than
291 initially suspected. We suggest that, in the absence of evidence for higher μ , fixing it and
292 estimating D may provide important clues regarding prior cryptic transmission. Considering
293 higher D values than suggested by case recognition is clearly relevant for epidemiological
294 investigations of outbreaks, as it widens the considered time window and may lead to identify
295 initially unsuspected sources of contamination.

296 The analysis of the 16 published outbreaks led to the definition of genetic thresholds that were
297 largely consistent with epidemiological evidence. For outbreaks 4 and 11, groupings were
298 discordant, as a lower threshold than initially used was inferred by our model. However, when
299 estimating the duration or substitution rate with our framework, higher values were observed
300 for both outbreaks, thus leading to group samples consistently with epidemiological evidence.
301 Outbreak 11 involved foodborne listeriosis with contaminated food where the two food
302 samples differed by 9 SNPs from the human samples, themselves separated by 2 SNPs. The
303 two food samples were isolated from two food outlets that had the same meat producer.
304 Because the incubation period of listeriosis is between 3 and 70 days and because intermittent
305 *L. monocytogenes* contamination during the production was observed [31], the duration of
306 contamination D might have been higher than initially defined, suggesting that the true
307 common ancestor of food and human isolates was in fact older than initially estimated from
308 the original publication. This illustrates the value of our estimation framework to inform

309 epidemiological investigations. Interestingly, when using model-estimated duration of
310 outbreak or substitution rate, we often observed an improved fit of the pairwise distance
311 distributions (**Table S2**).

312 For outbreak 8, low sequence data quality was observed for three genomes [24], including the
313 two genomes excluded from the outbreak by our model. Low quality data may have
314 artificially inflated their genetic distinctness, which underlines the importance of input
315 sequence data quality.

316 It is important to highlight the following limitations. First, all presented results were generated
317 by initialising the models with a fully homogeneous ancestral population. However, the
318 contaminating population may be slightly heterogeneous if it has a non-negligible population
319 size and had itself already evolved previously. In these cases, D might be interpreted as
320 incorporating the diversification time before source contamination. Second, we only modelled
321 mutation, neglecting other evolutionary processes such as recombination. Detection of
322 recombination among very closely related isolates is very challenging and its impact would be
323 limited. However, recombination with genetically distinct co-contaminants might occur and
324 recombined chromosomal regions should be removed from the analysis, especially when
325 using SNP-based analyses (by design, MLST moderates the impact of homologous
326 recombination). Third, the model does not incorporate demographic events within the
327 contaminated source, including population bottlenecks, which are potentially common in food
328 processing chains but which would be challenging to infer and to model. Finally, the
329 framework is designed for a single evolving population derived from a single bacterial
330 ancestor. When there is more than one contaminating genotype, our framework could be used
331 separately for each of these.

332 **Conclusions**

333 We describe an innovative approach to the ‘single strain’ definition using pathogen genomic
334 data by considering the most relevant features of specific outbreaks to define a credible
335 genetic distance threshold. This definition is grounded in evolutionary biology and alleviates
336 the need for *a priori* defined thresholds, which are not justified theoretically and may be
337 inappropriate in most cases. The inferred outbreak-tailored genetic thresholds provide a
338 reliable, non-arbitrary way to define epidemiologically related infection cases and to exclude
339 non-related sporadic strains. This approach is fast and easy to use. The additional ability to
340 estimate outbreak duration should also prove useful for point-source disease outbreak studies,

341 by providing a credible temporal window for epidemiological investigations aiming at
342 identifying and eliminating the sources of contaminations.

343 **Figure legends**

344

345 **Figure 1. Description of the framework.**

346 **A.** Left, threshold computation inputs: genetic distance matrix M , duration of outbreak D , set
347 of sample dates S_d , number of substitutions per site per year μ , and sequence length L (if based
348 on nucleotide sites) or number of genes g (if based on a gene-by-gene approach). Right,
349 model-based simulation: the algorithm is initialized with a homogenous population of
350 individuals. At each time step, substitutions are drawn from a Poisson distribution, until D is
351 reached. Samples are drawn randomly at the different observed sampling dates. A genetic
352 threshold is defined using *e.g.*, the 99th percentile of the distribution, and clusters of isolates
353 are derived by single linkage clustering, leading to rule-out non-outbreak isolates.

354 **B.** Left, the same model is used to estimate D or μ using MCMC, based on the following
355 inputs: the genetic distance matrix; the sampling dates; the sequence length; and either μ or D
356 (depending on which one is estimated).

357

358 **Figure 2. Assessment of the model's ability to classify outbreak isolates from the**
359 **simulation study.**

360 Specificity (**A**) and sensitivity (**B**) of isolates classification when $\mu = 8\text{e-}08$ (top) or $8\text{e-}07$
361 (bottom) substitutions per site per year. Each point provides specificity or sensitivity computed
362 from 20 independent outbreaks simulated with the same input parameters, with D ranging from
363 50 to 1000 days (x-axis) and R_d (the ratio of evolution duration between non-outbreak and
364 outbreak genomes) varying between 4.5 and 150 (colours). (**C**) 95% specificity threshold value
365 of D as a function of R_d (x-axis), computed for 9 values of μ (colours).

366

367 **Figure 3. Assessment of the quality of the estimation of D and μ through simulation.**

368 Precision of the estimation of D (**A**) and μ (**B**) from simulated data generated using different
369 values of D and μ . The sample size, defined as the number of observed samples and
370 associated dates, was set to $0.2xD$. On each panel, the upper banner and red line indicate the
371 expected D (**A**) and μ (**B**) values used to generate the simulated data. For each of the 240
372 synthetic outbreaks analysed, three independent MCMC chains were run, and the three

373 corresponding best estimates are shown (points). Vertical bars represent the average values of
374 the minimum and maximum of the 95% credible interval of the 3 MCMC chains. Each colour
375 corresponds to distinct values of μ or D used in the simulations (see keys).

376

377 **Figure 4. Impact of sampling density on the precision of the estimation of D and μ .**

378 The position of each symbol represents the difference between the expected and best estimates of D
379 (A) and μ (B) for each of 2400 outbreaks simulated using combinations of 4 values of D (60, 100, 200
380 and 400 days; represented in rows) and 3 values of μ (2E-07: blue diamonds, 4E-07: orange circles,
381 6E-07, red triangles; values in substitutions per site per year). Sampling density represents the
382 percentage of individuals sampled at each time step.

383

384 **Figure 5. Distance threshold derived from the modelling framework, and its effect on clustering:
385 example of outbreak 11.**

386 Panels A and C show the cgMLST distance distributions: observed distribution (blue, panels A and C),
387 simulated distribution without estimation (Orange, panel A), and simulated distribution using the
388 estimated duration of outbreak (red, panel C). Error bars represent the interval of prediction at 95% of
389 100 simulations. Blue vertical lines correspond to the derived distance threshold defined here as the
390 99th percentile of the distributions (A: from the observed distribution; C: from the simulated
391 distribution using the estimated duration of outbreak). Panels B and D show the single-linkage clusters
392 resulting from the derived distance threshold corresponding to panels A and C, respectively.

Table 1. Analysis of 16 outbreaks from literature. Genome length L is given in base pairs (bp), for outbreaks 5 and 8 to 16 the genome length was given by the number of loci g multiply by the average gene length. D_{lit} correspond to the duration of outbreak in days deduced from the published articles and μ_{lit} is the number of mutations per site per year found in the published article related or found in the literature in general. For each μ_{lit} value, the reference used is shown. $D_{estimated}$ (in days) and $\mu_{estimated}$ were estimated based on 3 MCMC chains; the associated 95% HPD for all the outbreaks as well as the corresponding genetic threshold is shown. Each threshold results from the 99th percentile of a pairwise differences distribution from 100 outbreak simulations.

Outbreak	Period	Country	Bacteria	Samples size		Source	Genomic marker	Ref. outbreak	L (bp)	D_{lit}	μ_{lit}	Ref. μ_{lit}	$D_{estimated}$	$\mu_{estimated}$	Cut-off		
				Human/Animal	Food										Using D_{lit} and μ_{lit}	Using $D_{estimated}$	Using $\mu_{estimated}$
1	Nov. 2016	Australia	SM^2	13	0	Chocolate mousse	SNP	[6]	4857450	120	12E10-7	[6]	95.64 (68.73-120.41)	1.01e-06 (4.83e-07-1.69e-06)	8	7	7
2	Jan-May 2014	Australia	SM^2	25	0	Chicken liver pâté	SNP	[22]	4857450	20	12E10-7	[6]	20.51 (20.01-54.7)	6.33e-07 (2.18e-07-1.09e-06)	1	1	1
3	Jan-May 2014	Australia	SM^2	20	0	Hot bread shop	SNP	[22]	4857450	38	12E10-7	[6]	44.92 (38.8-54.83)	1.13e-06 (2.78e-07-1.2e-06)	2	2	2
4	Oct-Nov 2013	Australia	CJ^3	7	2	Chicken liver pâté	SNP	[20]	1343000	9	3.23E10-5	[32]	28.08 (16.19-29.99)	1.38e-04 (9.11e-05-0.000312)	4	11	11
5	Dec 2002 - Jan 2003	Finland	CJ^3	2/2	2	Milk	cgMLST	[21]	1432000	65	3.23E10-5	[32]	73.07 (66.26-95.73)	3.53e-05 (2.19e-05-3.07e-04)	16	19	18
6	May-July 2011	Germany & France	EC^4 <i>O104:H4</i>	15	0	Sprouts	SNP	[19]	5437407	55	2.5E10-6	[33]	107.41 (70.33-180.02)	5.39e-06 (1.72e-06-5.90e-06)	7	13	13
7	2011	UK	EC^4 <i>O157</i>	10	0	Unwashed	SNP	[23]	4122236	340	2.26E10-7	[34]	352.	3.21e-07	2	2	2

	vegetables											(340.07-394.00)	(1.16e-07-5.26e-07)				
8	2012-2013	B ¹	LM ⁵	5	10	Beef	cgMLST	[24]	1462000	466	4.3E10-7	[35]	494.49 (466.09-582.86)	2.09e-06 (3.44e-07-3.78e-06)	2	2	7
9	2007-2013	B ¹	LM ⁵	5	3	Crabmeat	cgMLST	[24]	1732000	2200	4.3E10-7	[35]	2270.8 (2200.54-2569.1)	4.83e-07 (2.75e-07-9.93e-07)	12	12	13
10	2013-2014	B ¹	LM ⁵	5	4	Sandwiches	cgMLST	[24]	1698000	289	4.3E10-7	[35]	473.51 (318.76-499.83)	1.05e-06 (3.78e-07-3.63e-06)	2	4	5
11	2013-2014	B ¹	LM ⁵	2	2	Ox tongue	cgMLST	[24]	1464000	943	4.3E10-7	[35]	1559.36 (1025.81-1599.92)	1.17e-06 (6.22e-07-4.10e-06)	4	7	10
12	2009-2011	B ¹	LM ⁵	9	1	Unknown	cgMLST	[24]	1685000	783	4.3E10-7	[35]	805.65 (783.08-888.25)	2.46e-07 (1.4e-07-4.52e-07)	6	6	4
13	2013	T ¹	LM ⁵	4	1	Rakfisk	cgMLST	[24]	1526000	180	4.3E10-7	[35]	186.32 (76.94-299.33)	5.47e-07 (8.25e-08-2.98e-06)	1	1	1
14	2013-2014	X ¹	LM ⁵	13	6	Foie gras	cgMLST	[24]	1686000	161	4.3E10-7	[35]	492.69 (234.34-499.94)	2.91e-06 (1.1e-06-4.29e-06)	2	5	8
15	2012	X ¹	LM ⁵	4	9	Cheese	cgMLST	[24]	1698000	548	4.3E10-7	[35]	384.79 (244.07-558.99)	3.54e-07 (1.53e-07-5.35e-07)	5	4	4
16	2012	C ¹	LM ⁵	25	0	Brie cheese	cgMLST	[24]	1707000	150	4.3E10-7	[35]	294.28 (252.2-299.99)	3.43e-06 (1.81e-06-4.3e-06)	2	3	9

¹ Country code from the reference article; ² SM: *Salmonella enterica* serovar *Typhimurium*; ³ CJ: *Campylobacter jejuni*; ⁴ EC: *Escherichia coli*; ⁵ LM: *Listeria monocytogenes*

Acknowledgements

We acknowledge the financial support of the MedVetKlebs project, a component of European Joint Programme One Health EP, which has received funding from the European Union Horizon 2020 Research and Innovation Programme (Grant N. 773830). We thank Xavier Didelot and Olivier Tenaillon for critical feedback on an earlier version of the manuscript, and Chiara Crestani for help in manuscript and figures formatting.

Authors contributions

S.B. conceived the presented approach. A.D., L.O. and S.B. designed the model and the computational framework. A.D. programmed the model and analysed data and simulations. A.D., L.O. and S.B. interpreted the results. A.D., L.O. and S.B. wrote the manuscript. The three authors read and approved the final manuscript.

Declaration of interests

All authors report no competing interests.

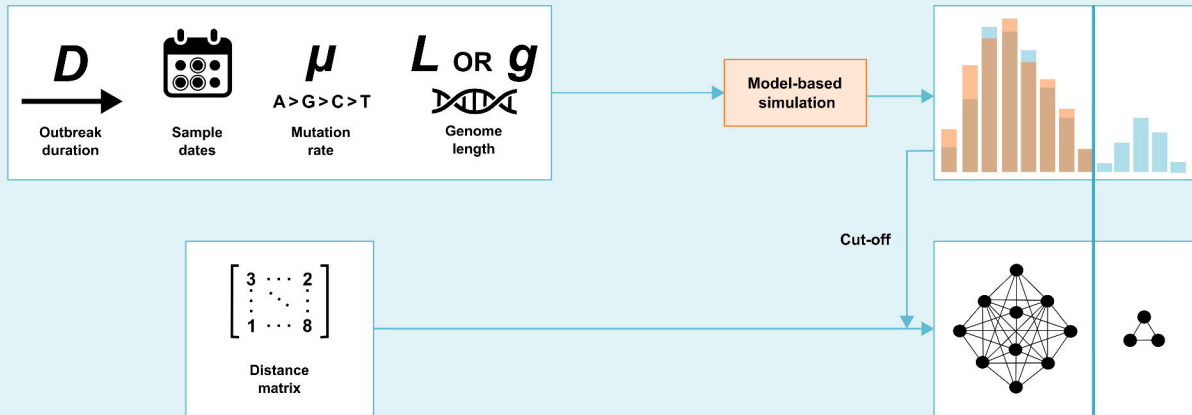
References

1. Gerner-Smidt P, Hise K, Kincaid J, Hunter S, Rolando S, Hyttiä-Trees E, et al. PulseNet USA: a five-year update. *Foodborne Pathog Dis.* 2006;3: 9–19. doi:10.1089/fpd.2006.3.9
2. van Belkum A, Tassios PT, Dijkshoorn L, Haeggman S, Cookson B, Fry NK, et al. Guidelines for the validation and application of typing methods for use in bacterial epidemiology. *Clin Microbiol Infect.* 2007;13 Suppl 3: 1–46.
3. Ruan Z, Yu Y, Feng Y. The global dissemination of bacterial infections necessitates the study of reverse genomic epidemiology. *Brief Bioinform.* 2020;21: 741–750. doi:10.1093/bib/bbz010
4. Tenover FC, Arbeit RD, Goering RV, Mickelsen PA, Murray BE, Persing DH, et al. Interpreting chromosomal DNA restriction patterns produced by pulsed-field gel electrophoresis: criteria for bacterial strain typing. *J Clin Microbiol.* 1995;33: 2233–9.
5. Pightling AW, Pettengill JB, Luo Y, Baugher JD, Rand H, Strain E. Interpreting Whole-Genome Sequence Analyses of Foodborne Bacteria for Regulatory Applications and Outbreak Investigations. *Front Microbiol.* 2018;9. doi:10.3389/fmicb.2018.01482
6. Octavia S, Wang Q, Tanaka MM, Kaur S, Sintchenko V, Lan R. Delineating Community Outbreaks of *Salmonella enterica* Serovar Typhimurium by Use of Whole-Genome Sequencing: Insights into Genomic Variability within an Outbreak. *J Clin Microbiol.* 2015;53: 1063–1071. doi:10.1128/JCM.03235-14
7. Schürch AC, Arredondo-Alonso S, Willems RJL, Goering RV. Whole genome sequencing options for bacterial strain typing and epidemiologic analysis based on single nucleotide polymorphism versus gene-by-gene-based approaches. *Clin Microbiol Infect.* 2018;24: 350–354. doi:10.1016/j.cmi.2017.12.016
8. Schürch AC, Siezen RJ. Genomic tracing of epidemics and disease outbreaks. *Microb Biotechnol.* 2010;3: 628–633. doi:10.1111/j.1751-7915.2010.00224.x
9. Reuter S, Ellington MJ, Cartwright EJP, Köser CU, Török ME, Gouliouris T, et al. Rapid bacterial whole-genome sequencing to enhance diagnostic and public health microbiology. *JAMA Intern Med.* 2013;173: 1397–1404. doi:10.1001/jamainternmed.2013.7734
10. Moura A, Tourdjman M, Leclercq A, Hamelin E, Laurent E, Fredriksen N, et al. Real-Time Whole-Genome Sequencing for Surveillance of *Listeria monocytogenes*, France. *Emerg Infect Dis.* 2017;23: 1462–1470. doi:10.3201/eid2309.170336
11. van Belkum A, Tassios PT, Dijkshoorn L, Haeggman S, Cookson B, Fry NK, et al. Guidelines for the validation and application of typing methods for use in bacterial epidemiology. *Clin Microbiol Infect.* 2007;13 Suppl 3: 1–46. doi:10.1111/j.1469-0691.2007.01786.x
12. Schürch AC, Arredondo-Alonso S, Willems RJL, Goering RV. Whole genome sequencing options for bacterial strain typing and epidemiologic analysis based on single nucleotide polymorphism versus gene-by-gene-based approaches. *Clin Microbiol Infect.* 2018;24: 350–354. doi:10.1016/j.cmi.2017.12.016
13. Collineau L, Boerlin P, Carson CA, Chapman B, Fazil A, Hetman B, et al. Integrating Whole-Genome Sequencing Data Into Quantitative Risk Assessment of Foodborne Antimicrobial Resistance: A Review of Opportunities and Challenges. *Front Microbiol.* 2019;10: 1107. doi:10.3389/fmicb.2019.01107

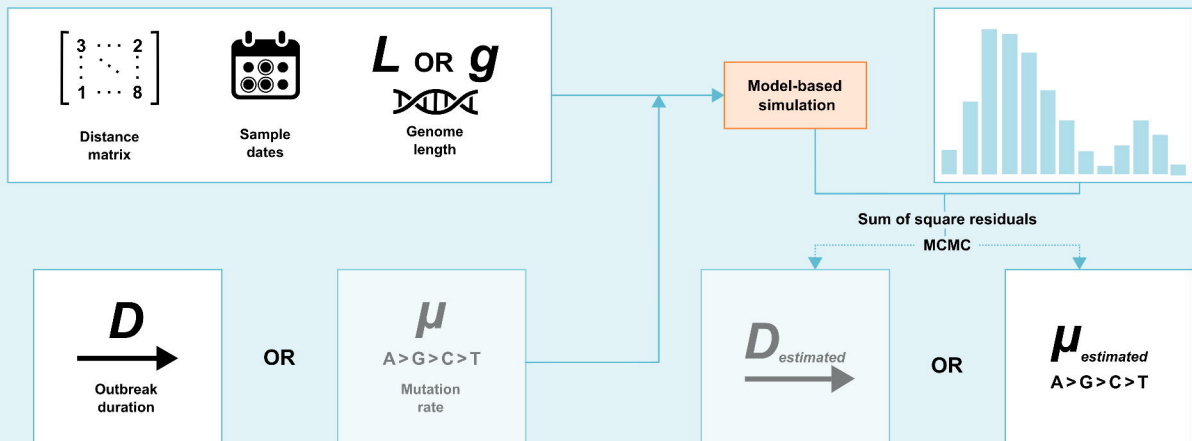
14. Stimson J, Gardy J, Mathema B, Crudu V, Cohen T, Colijn C. Beyond the SNP Threshold: Identifying Outbreak Clusters Using Inferred Transmissions. *Mol Biol Evol.* 2019;36: 587–603. doi:10.1093/molbev/msy242
15. Gordon NC, Pichon B, Golubchik T, Wilson DJ, Paul J, Blanc DS, et al. Whole-Genome Sequencing Reveals the Contribution of Long-Term Carriers in *Staphylococcus aureus* Outbreak Investigation. *J Clin Microbiol.* 2017;55: 2188–2197. doi:10.1128/JCM.00363-17
16. Coll F, Raven KE, Knight GM, Blane B, Harrison EM, Leek D, et al. Definition of a genetic relatedness cutoff to exclude recent transmission of meticillin-resistant *Staphylococcus aureus*: a genomic epidemiology analysis. *Lancet Microbe.* 2020;1: e328–e335. doi:10.1016/S2666-5247(20)30149-X
17. Wright S. Evolution in Mendelian Populations. *Genetics.* 1931;16: 97–159.
18. Fisher RA. XXI.—On the Dominance Ratio. *Proceedings of the Royal Society of Edinburgh.* 1923;42: 321–341. doi:10.1017/S0370164600023993
19. Grad YH, Lipsitch M, Feldgarden M, Arachchi HM, Cerqueira GC, Fitzgerald M, et al. Genomic epidemiology of the *Escherichia coli* O104:H4 outbreaks in Europe, 2011. *Proc Natl Acad Sci USA.* 2012;109: 3065–3070. doi:10.1073/pnas.1121491109
20. Moffatt CRM, Greig A, Valcanis M, Gao W, Seemann T, Howden BP, et al. A large outbreak of *Campylobacter jejuni* infection in a university college caused by chicken liver pâté, Australia, 2013. *Epidemiology & Infection.* 2016;144: 2971–2978. doi:10.1017/S0950268816001187
21. Revez J, Zhang J, Schott T, Kivistö R, Rossi M, Hänninen M-L. Genomic variation between *Campylobacter jejuni* isolates associated with milk-borne-disease outbreaks. *J Clin Microbiol.* 2014;52: 2782–2786. doi:10.1128/JCM.00931-14
22. Phillips A, Sotomayor C, Wang Q, Holmes N, Furlong C, Ward K, et al. Whole genome sequencing of *Salmonella* Typhimurium illuminates distinct outbreaks caused by an endemic multi-locus variable number tandem repeat analysis type in Australia, 2014. *BMC Microbiol.* 2016;16: 211. doi:10.1186/s12866-016-0831-3
23. Holmes A, Allison L, Ward M, Dallman TJ, Clark R, Fawkes A, et al. Utility of Whole-Genome Sequencing of *Escherichia coli* O157 for Outbreak Detection and Epidemiological Surveillance. *J Clin Microbiol.* 2015;53: 3565–3573. doi:10.1128/JCM.01066-15
24. Nielsen EM, Björkman JT, Kiil K, Grant K, Dallman T, Painset A, et al. Closing gaps for performing a risk assessment on *Listeria monocytogenes* in ready-to-eat (RTE) foods: activity 3, the comparison of isolates from different compartments along the food chain, and from humans using whole genome sequencing (WGS) analysis. *EFSA Supporting Publications.* 2017;14: 1151E. doi:10.2903/sp.efsa.2017.EN-1151
25. Laughlin M, Bottichio L, Weiss J, Higa J, McDonald E, Sowadsky R, et al. Multistate outbreak of *Salmonella* Poona infections associated with imported cucumbers, 2015–2016. *Epidemiol Infect.* 2019;147. doi:10.1017/S0950268819001596
26. McCollum JT, Cronquist AB, Silk BJ, Jackson KA, O'Connor KA, Cosgrove S, et al. Multistate Outbreak of Listeriosis Associated with Cantaloupe. *New England Journal of Medicine.* 2013;369: 944–953. doi:10.1056/NEJMoa1215837
27. Multi-country outbreak of *Salmonella* Stanley infections – Third update. *EFSA Supporting Publications.* 2014;11: 592E. doi:10.2903/sp.efsa.2014.EN-592

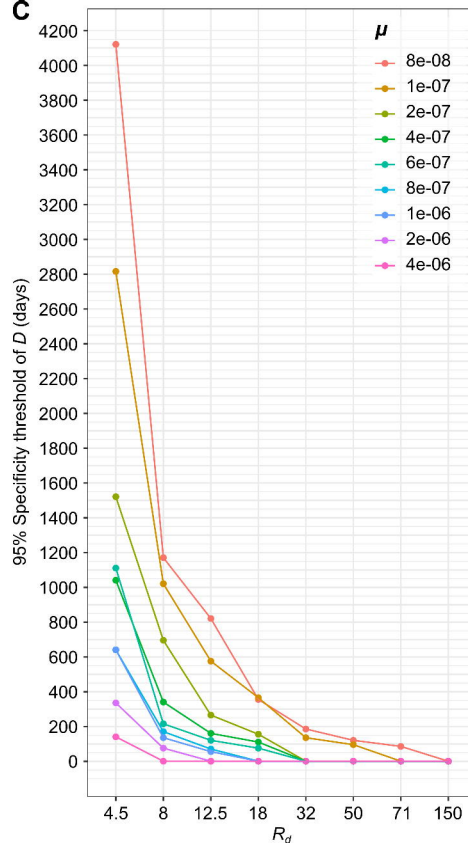
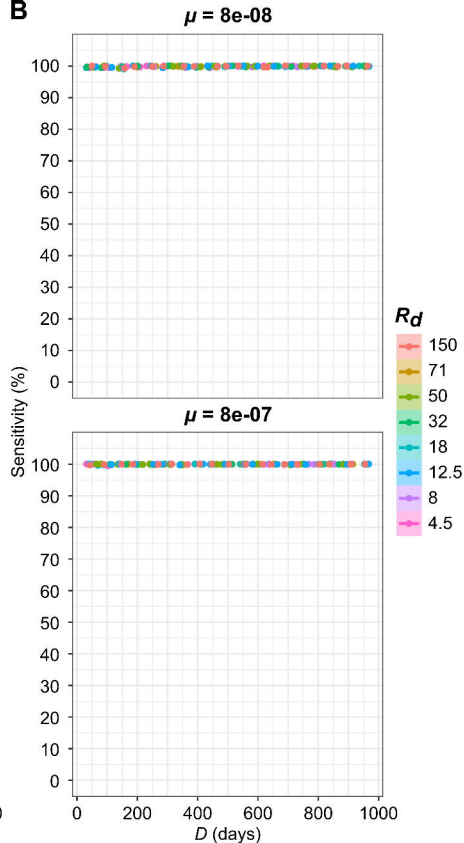
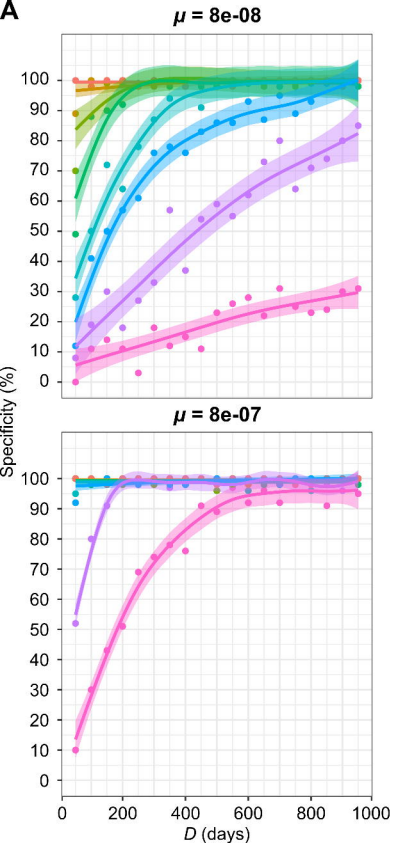
28. Coll F, Raven KE, Knight GM, Blane B, Harrison EM, Leek D, et al. Definition of a genetic relatedness cutoff to exclude recent transmission of meticillin-resistant *Staphylococcus aureus*: a genomic epidemiology analysis. *The Lancet Microbe*. 2020;1: e328–e335. doi:10.1016/S2666-5247(20)30149-X
29. Perrin A, Larssonneur E, Nicholson AC, Edwards DJ, Gundlach KM, Whitney AM, et al. Evolutionary dynamics and genomic features of the *Elizabethkingia anophelis* 2015 to 2016 Wisconsin outbreak strain. *Nature Communications*. 2017;8: 15483. doi:10.1038/ncomms15483
30. Duchêne S, Holt KE, Weill F-X, Le Hello S, Hawkey J, Edwards DJ, et al. Genome-scale rates of evolutionary change in bacteria. *Microb Genom*. 2016;2. doi:10.1099/mgen.0.000094
31. Lamden KH, Fox AJ, Amar CFL, Little CL. A case of foodborne listeriosis linked to a contaminated food-production process. *Journal of Medical Microbiology*. 2013;62: 1614–1616. doi:10.1099/jmm.0.064055-0
32. Wilson DJ, Gabriel E, Leatherbarrow AJH, Cheesbrough J, Gee S, Bolton E, et al. Rapid Evolution and the Importance of Recombination to the Gastroenteric Pathogen *Campylobacter jejuni*. *Mol Biol Evol*. 2009;26: 385–397. doi:10.1093/molbev/msn264
33. Grad YH, Godfrey P, Cerquiera GC, Mariani-Kurkdjian P, Gouali M, Bingen E, et al. Comparative Genomics of Recent Shiga Toxin-Producing *Escherichia coli* O104:H4: Short-Term Evolution of an Emerging Pathogen. *mBio*. 2013;4. doi:10.1128/mBio.00452-12
34. Reeves PR, Liu B, Zhou Z, Li D, Guo D, Ren Y, et al. Rates of Mutation and Host Transmission for an *Escherichia coli* Clone over 3 Years. *PLoS One*. 2011;6. doi:10.1371/journal.pone.0026907
35. Halbedel S, Prager R, Fuchs S, Trost E, Werner G, Flieger A. Whole-Genome Sequencing of Recent *Listeria monocytogenes* Isolates from Germany Reveals Population Structure and Disease Clusters. *J Clin Microbiol*. 2018;56. doi:10.1128/JCM.00119-18

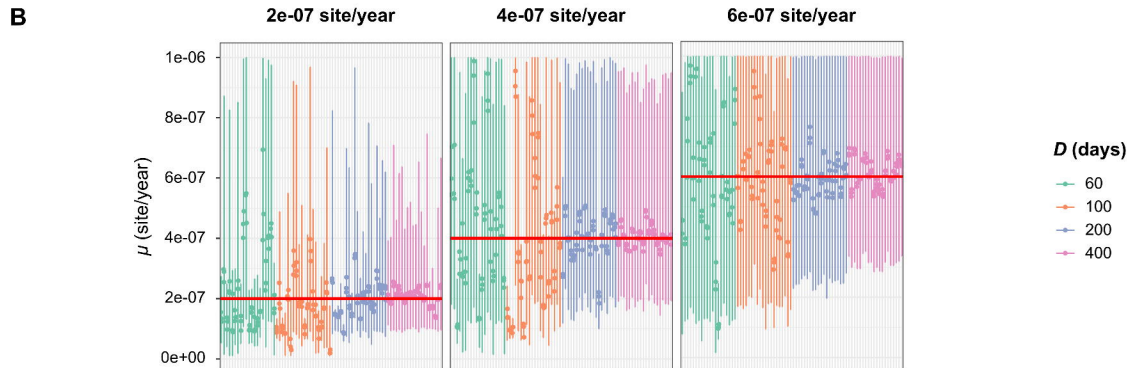
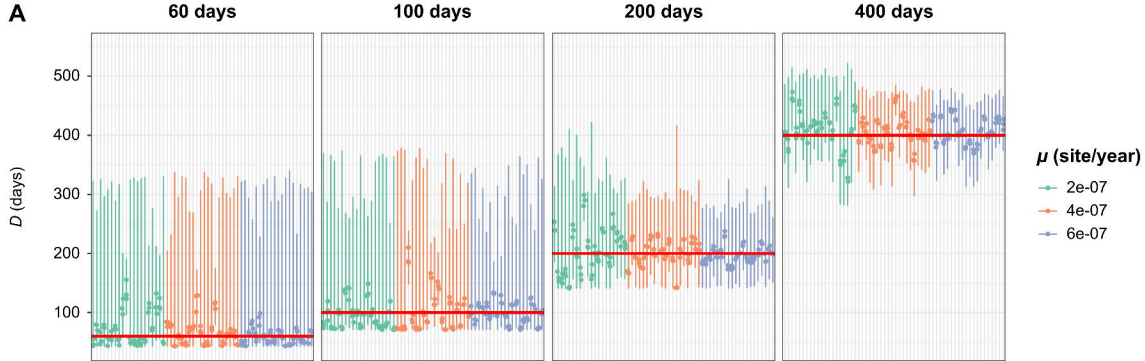
A

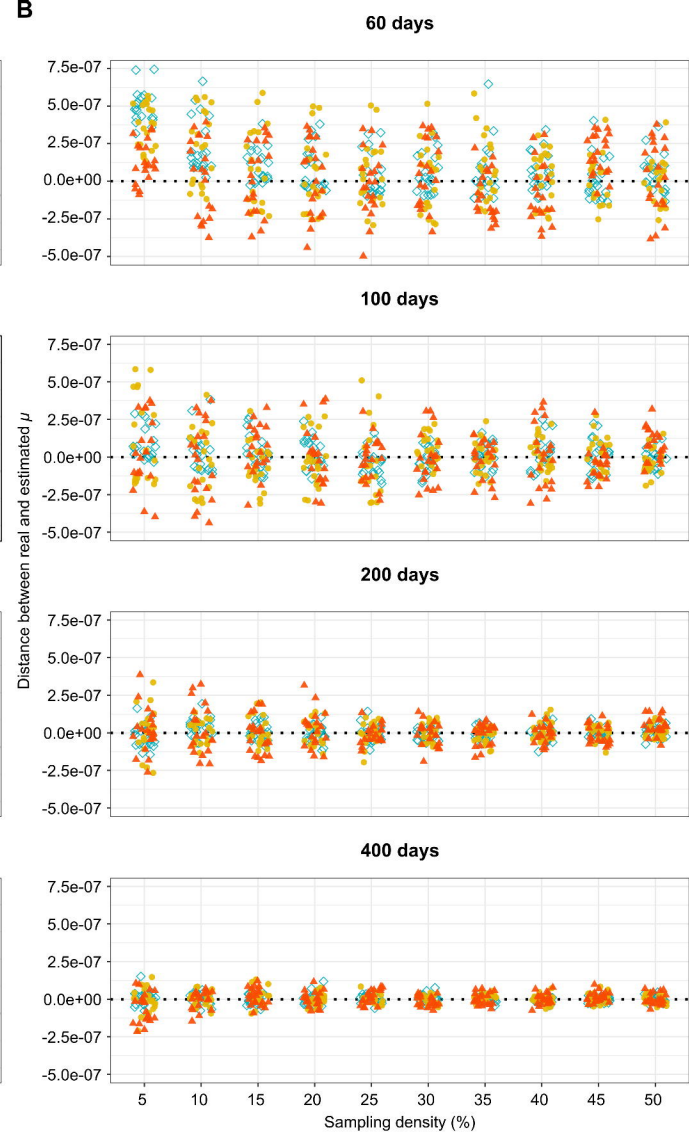
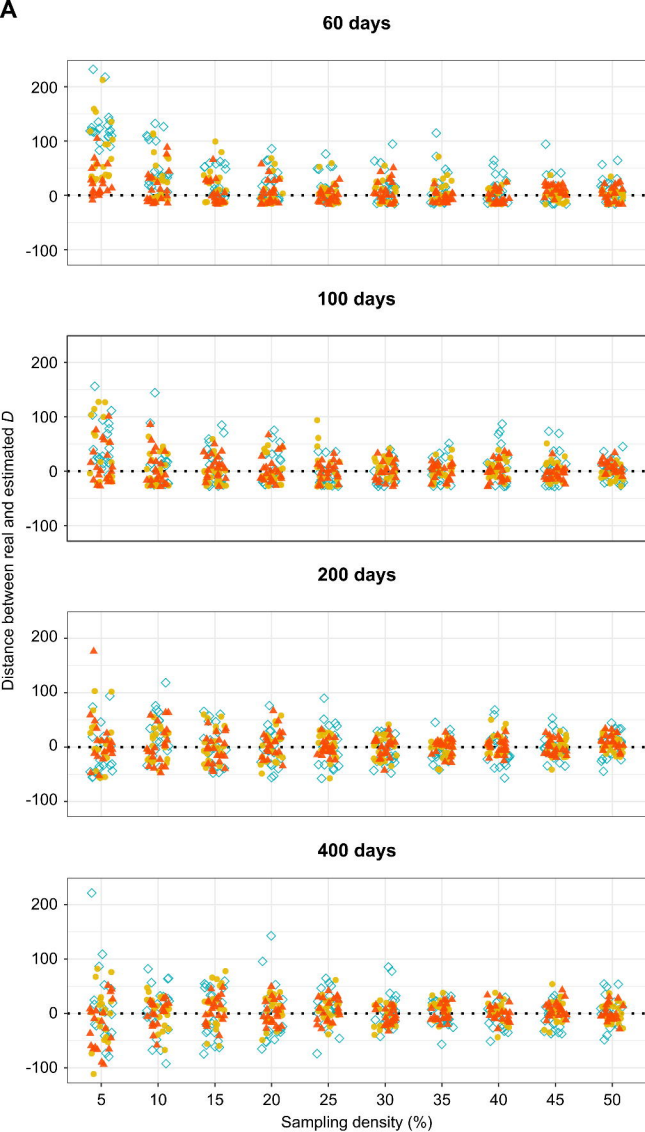


B









μ (site/year)

◇ 2e-07

● 4e-07

▲ 6e-07

