

1 **Short title:**

2 Bayesian ancestral state reconstruction models for investigating *Salmonella* outbreaks

3

4 **Long title:**

5 Investigation of the validity of two Bayesian ancestral state reconstruction models for
6 estimating *Salmonella* transmission during outbreaks

7

8 **Authors:**

9 Samuel J. Bloomfield

10 Quadram Institute, Norwich Research Park, Colney Lane, Norwich, United Kingdom

11 Contribution: Conceptualization, Data collection, Formal analysis, Investigation,

12 Methodology, Project Administration, Software, Writing – Original Draft preparation,

13 Writing – Review and Editing,

14

15 Timothy G. Vaughan

16 Department of Biosystems Science and Engineering, ETH Zurich, Zurich, Switzerland

17 Contribution: Conceptualization, Data collection, Formal analysis, Investigation,

18 Methodology, Software, Supervision, Writing – Review and Editing,

19

20 Jackie Benschop

21 Molecular Epidemiology and Public Health Laboratory, Massey University, Palmerston

22 North, New Zealand

23 Contribution: Funding acquisition, Supervision, Writing – Review and Editing,

24

25 Jonathan C. Marshall

26 Molecular Epidemiology and Public Health Laboratory, Massey University, Palmerston
27 North, New Zealand

28 Contribution: Data collection, Formal analysis, Methodology, Supervision, Writing – Review
29 and Editing,
30

31 David T. S. Hayman

32 Molecular Epidemiology and Public Health Laboratory, Massey University, Palmerston
33 North, New Zealand

34 Contribution: Conceptualization, Data collection, Formal analysis, Methodology,
35 Supervision, Writing – Review and Editing,
36

37 Patrick J. Biggs

38 Molecular Epidemiology and Public Health Laboratory, Massey University, Palmerston
39 North, New Zealand

40 Contribution: Supervision, Visualization, Writing – Review and Editing,
41

42 Philip E. Carter

43 Institute of Environmental Science and Research, Keneperu, New Zealand

44 Contribution: Supervision, Writing – Review and Editing,
45

46 Nigel P. French

47 Molecular Epidemiology and Public Health Laboratory, Massey University, Palmerston
48 North, New Zealand. The New Zealand Foo Safety and Science Centre.

49 Contribution: Conceptualization, Data collection, Formal analysis, Funding acquisition,
50 Methodology, Project Administration, Supervision, Writing – Review and Editing,

51 **Abstract**

52 Ancestral state reconstruction models use genetic data to characterize a group of
53 organisms' common ancestor. These models have been applied to salmonellosis outbreaks to
54 estimate the number of transmissions between different animal species that share similar
55 geographical locations, with animal host as the state. However, as far as we are aware, no
56 studies have validated these models for outbreak analysis. In this study, salmonellosis
57 outbreaks were simulated using a stochastic Susceptible-Infected-Recovered model, and the
58 host population and transmission parameters of these simulated outbreaks were estimated
59 using Bayesian ancestral state reconstruction models (discrete trait analysis (DTA) and
60 structured coalescent (SC)). These models were unable to accurately estimate the number of
61 transmissions between the host populations or the amount of time spent in each host
62 population. The DTA model was inaccurate because it assumed the number of isolates
63 sampled from each host population was proportional to the number of individuals infected
64 within each host population. The SC model was inaccurate possibly because it assumed that
65 each host population's effective population size was constant over the course of the simulated
66 outbreaks. This study highlights the need for phylodynamic models that can take into
67 consideration factors that influence the characteristics and behavior of outbreaks, e.g.
68 changing effective population sizes, variation in infectious periods, intra-population
69 transmissions, and disproportionate sampling of infected individuals.

70

71 **Introduction**

72 Ancestral state reconstruction models estimate the ancestral states of organisms based
73 on their evolutionary history. Outbreaks are "...the occurrence of disease in excess of what
74 would normally be expected in a defined community, geographical area or season" (1).
75 Ancestral state reconstruction models have been used to investigate the transmission of

76 infectious agents between animal populations over the course of outbreaks, with host
77 population as the state (2). However, as far as we are aware, no studies have validated these
78 models for this type of analysis.

79 The discrete trait analysis (DTA) and structured coalescent (SC) models are ancestral
80 state reconstruction models. Both models treat each host population as a discrete trait and can
81 be approximated using Markov chain Monte Carlo methods (3,4). There are many differences
82 between these two ancestral state reconstruction models. In the context of host association
83 studies, the DTA model uses a substitution model to model the transmission between host
84 populations (3). The pruning algorithm (5), often used in phylogenetic analysis to account for
85 possible mutations, is similarly used by the DTA model to integrate all possible migration
86 histories (6). The SC model assumes that the pathogen associated with each host population
87 has a fixed effective population size and models the transmission between populations. The
88 DTA model assumes that the number of offspring an individual pathogen is likely to produce
89 is independent of its host population, whilst the SC model allows for variation between host
90 populations (4). The DTA model assumes that the proportion of isolates sampled from each
91 host population is proportional to the size of the pathogen population associated with that
92 host, whilst the SC model allows for variation in these population sizes (6). Some of these
93 assumptions are applicable to the investigation of outbreaks (e.g. varying effective population
94 size), whilst others are not (e.g. isolate proportionality).

95 Salmonellosis is an intestinal infection caused by non-typhoidal *Salmonella* strains.
96 Salmonellosis outbreaks vary in size and can involve one or more host populations (7).
97 Identifying the amount of time *Salmonella* spends in a host population over an outbreak and
98 the amount of transmission between host populations can inform control measures to limit
99 salmonellosis outbreaks, e.g. if human cases are primarily from exposure to poultry sources
100 then control measures that limit human exposure to poultry or decrease the amount of

101 *Salmonella* in poultry may be beneficial. However, there is growing evidence that exposure
102 to human sources contributes more to salmonellosis outbreaks than previously thought (8).
103 Therefore, methods and models are required that can approximate the number of cases that
104 are the result of exposure to different animal and/or human sources. The aim of this study
105 was to use simulated outbreaks to investigate whether the DTA or SC models could be
106 applied to infer transmission dynamics in outbreaks involving multiple hosts, motivated by
107 non-typhoidal *Salmonella*.

108

109 **Methods**

110 **Outbreak simulations**

111 The MASTER package (9) in BEAST2 (10) was used to simulate stochastic
112 transmission dynamics for a pathogen infecting structured populations, including associated
113 phylogenetic and transmission trees. Outbreaks were generated using a stochastic
114 Susceptible-Infected-Recovered (SIR) model, intended to simulate the transmission of
115 zoonotic salmonellosis. In this model, susceptible host individuals become infectious by
116 exposure to other infected individuals:

$$117 \quad S_i + I_i \xrightarrow{\beta_{ii}} 2I_i \quad (1)$$

$$118 \quad S_i + I_j \xrightarrow{\beta_{ji}} I_i + I_j \quad (2)$$

119

120 Equation 1 represents the transmission of the infectious agent from an infected
121 individual to a susceptible individual of the same host population. Equation 2 represents the
122 transmission of the infectious agent from an infected individual to a susceptible individual of
123 another host population. Here, S_i represents a susceptible individual from one host
124 population, I_i represents an infectious individual from the same host population, I_j represents

125 an infectious individual from another host population, and β_{ii} and β_{ji} represents the
126 transmission rate per susceptible individual per infectious individual.

127 In this model, infectious individuals also recover or are removed over time:

$$128 \quad I_i \xrightarrow{\gamma_i} R_i \quad (3)$$

129 Equation 3 determines the infectious period for an infectious individual. Here, I_i
130 represents an infectious individual in one host population, R_i represents a recovered/removed
131 individual in the same host population, and γ_i represents the recovery/removal rate per
132 infectious individual for this host population. The mean infectious period for a host of type i
133 is $\frac{1}{\gamma_i}$.

134

135 **Simulated outbreaks**

136 We simulated 23 outbreaks using the MASTER package, hereinafter ‘outbreak
137 simulations’. This created 23 transmission trees consisting of all the transmissions that took
138 place over the course of each simulated outbreak (Fig 1). These simulations consisted of two
139 host populations: human and animal. We wanted to compare the simulated outbreaks with a
140 previously reported salmonellosis outbreak in New Zealand that involved *Salmonella*
141 *enterica* serovar Typhimurium DT160 (herein, DT160) (11). Therefore, the initial susceptible
142 host population size, infectious period (γ) and transmission rate (β) values varied between
143 the 23 simulations but represented possible values for salmonellosis outbreaks in New
144 Zealand (S1 Appendix).

145

146 **Fig 1.** Flow diagram of the methods used to compare the SC and DTA models using
147 various sampling methods. White rectangles represent the methods used and blue rectangles
148 represent the data produced.

149

150 **Simulated genetic sequences from outbreaks**

151 One hundred '*Salmonella*' isolates were randomly sampled from each outbreak
152 simulation, after stratifying for host population, hereinafter 'random sampling'. For each
153 outbreak simulation, the transmission tree was simplified to only include nodes common to
154 the 100 isolates (both steps were accomplished using custom Perl scripts). The sampled
155 transmission trees were used to simulate genetic data for the 23 simulated outbreaks using the
156 sequence simulation capability of the BEAST 2 package MASTER, hereinafter 'sequence
157 simulations'. 800 SNPs were simulated in total for the 100 isolates, similar to the 793 core
158 SNPs shared by 109 DT160 isolates (11). Perl and R scripts were used to analyze the sampled
159 transmission tree and to calculate the amount of time spent in each host population and
160 quantify the number of transmissions, later referred to as the 'known parameters'.

161

162 **Model consistency**

163 To investigate variation in model estimates between different samples (i.e. model
164 consistency), one of the simulated outbreaks was randomly sampled 10 times after stratifying
165 for host population. For each sample, sequence simulations were used to create genetic data.

166

167 **Sample size**

168 To investigate the effect of different sample sizes on the models' estimates, one of the
169 simulated outbreaks was randomly sampled 12 times. The number of isolates sampled
170 systematically ranged from 25 to 300 isolates in 11 increments of 25. For each sample,
171 sequence simulations were used to create genetic data. The genetic data systematically ranged
172 from 200 to 2400 SNPs in 11 increments of 200, respectively. To determine if sample size
173 affected the extremity of a model's estimates, the simulated outbreak chosen had significantly

174 different population values between host populations and similar transmission values for
175 comparison.

176

177 **Disproportionate sampling**

178 To investigate the effect of the relative number of isolates from each source on model
179 estimates (i.e. disproportionate sampling), as expected during the outbreaks, one of the
180 simulated outbreaks was randomly sampled 10 times with different numbers of animal and
181 human isolates. For each sample, 100 isolates were analyzed, but the proportion of isolates
182 that were from each host population were systematically ranged from 5-95% in 10%
183 intervals. For each sample, sequence simulations were used to create genetic data.

184

185 **Equal-time sampling**

186 To investigate an alternative sampling method, 'equal-time sampling', an in-house
187 Perl script was used to stratify the isolates from the initial 23 simulated outbreaks by host
188 population, before randomly sampling an equal number of isolates from each year of the
189 simulated outbreaks, to a total of 100 isolates. Sequence simulations were used to create
190 genetic data for the samples.

191

192 **Equal intra-population transmission and infectious periods**

193 To investigate if different intra-population transmission rates and infectious periods
194 had any effect on model estimates, twelve additional outbreaks were simulated but with equal
195 intra-population transmission rates and infectious periods (EPTI) for both host populations,
196 but inter-population transmission rates and initial susceptible host population sizes that
197 varied. For each simulation, 100 isolates were sampled using random sampling, and sequence
198 simulations were used to create genetic data.

199

200 **DTA model**

201 For the DTA model, the genetic data was imported into BEAUti 1.8.3 to create an
202 XML file for BEAST 1.8.3 (12). The generalized time reversible (GTR) model was used to
203 model base substitutions (13), the Gaussian Markov random field (GMRF) Bayesian skyride
204 model was used to allow for changes in the effective population size (14), and a strict
205 molecular clock was used to estimate the mutation rate, which was calibrated by the tip date.
206 The XML file was run in BEAST for 10 million steps as a single run with a 10% burn-in.

207

208 **SC model**

209 For the SC model, the genetic data was imported into BEAUti 2.4 with the
210 MultiTypeTree package (4) to create an XML file for BEAST 2.4 (10). The GTR model was
211 used to model base substitution and a strict molecular clock was used to estimate the
212 mutation rate, which was calibrated by the tip date. The XML file was run in BEAST for 250
213 million steps as a single run with a 10% burn-in. The SC model was run for a larger number
214 of steps than the DTA model as its population and transmission parameters took longer to
215 converge. BEAST 1.8.3 is unable to run the SC model, unlike BEAST 2.4. BEAST 2.4 can
216 run GMRF and DTA models but does not have a BEAUti interface to easily set up these
217 models. BEAST 1.8.3. does have an interface for these models so was used for the DTA
218 model.

219

220 **Model comparison**

221 The SC and DTA models were used to estimate the amount of time spent in each host
222 population (population parameters) and the amount of transmissions between the host
223 populations (transmission parameters). However, the models' raw outputs were not directly

224 comparable, as the SC model's implementation explicitly records transmissions along
225 branches, whilst the DTA approach integrates and marginalizes over these transmissions and
226 therefore does not record them in its output. Therefore, the relative amount of time (i.e.
227 proportion) spent in each host population and the relative number of inter-population
228 transmissions made up of each transmission were compared. The performance of the two
229 models were compared using four parameters:

- 230 1. The proportion of outbreak simulations that a model included the known parameter
231 within their 95% highest posterior density (HPD) intervals.
- 232 2. The mean squared error between a known parameter and a model's mean estimates.
- 233 3. The size of a model's 95% HPD intervals.
- 234 4. The correlation coefficient between a known parameter and a model's mean estimates.

235

236 **DT160 outbreak**

237 The DTA and SC models were used to analyze a previously-described salmonellosis
238 outbreak in New Zealand caused by DT160 (11). 109 DT160 isolates from animal (n=74) and
239 human (n=35) host populations over 14 years were investigated using the 793 core SNPs they
240 shared.

241

242 **Scripts**

243 The in-house scripts used in this study are available from GitHub
244 (<https://github.com/samuelbloomfield/Scripts-for-outbreak-simulations>).

245

246 **Results**

247 **Model consistency**

248 There was some variation in the DTA and SC models' population and transmission
249 mean estimates for the same simulated outbreak that was randomly sampled ten times (Fig 2).
250 The SC model's 95% HPD intervals included known population parameters more frequently,
251 whilst the DTA model's 95% HPD intervals included known transmission parameters more
252 frequently.

253 The outbreak transmission tree was the same for the ten samples, as these samples
254 were taken from the same simulated outbreak. However, the samples consisted of different
255 animal and human isolates, such that when the outbreak transmission tree was simplified to
256 only include nodes and branches common to these isolates, there was some variation in the
257 time spent in animal and human populations, and the number of transmissions between these
258 populations between samples. The known parameters were taken from the ten sampled
259 transmission trees, not the entire outbreak transmission tree, resulting in slight differences in
260 the known parameters between the ten samples. This is true for other analyses below that
261 sampled the same outbreak multiple times. Some of the outbreaks investigated in this
262 outbreak consisted of hundreds of thousands of infected animals and humans (S1 Appendix),
263 leaving large outbreak transmission trees that required large time periods to calculate the
264 number of transmissions and time spend in the populations. The small amount of variation in
265 the sampled transmission trees and the outbreak transmission tree for this dataset suggests
266 that the sampled transmission tree parameters are representative of the outbreak transmission
267 tree parameters.

268

269 **Fig 2.** The proportion of time spent in the animal (A and E) and human (B and F) host
270 populations, and the proportion of inter-population transmissions made up of animal-to-
271 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
272 (blue: A-D) and DTA (red: E-F) models, for 10 random samples of the same simulated

273 outbreak. The circles represent the mean, the error bars represent the 95% HPD interval, the
274 black horizontal lines represent the known parameters for the sampled outbreaks, and the
275 grey horizontal lines represent the known parameters for the entire outbreak.

276

277 **Sample size**

278 The DTA and SC models were affected by variation in sample size for the same
279 simulated outbreak differently. Increased sample sizes were associated with smaller 95%
280 HPD intervals and more accurate and extreme mean population estimates by the SC model up
281 to 100 samples. After this point, increased sample sizes had little effect on the precision,
282 extremity or accuracy of the model's mean population estimates (Fig 3). The DTA model's
283 mean population estimates were more precise than the SC model's. Sample size had no effect
284 on their accuracy but decreased the size of their 95% HPD intervals. The accuracy of the SC
285 and DTA models' mean transmission estimates and their 95% HPD intervals displayed some
286 variation, but there were no trends with sample size.

287

288 **Fig 3.** The proportion of time spent in the animal (A and E) and human (B and F) host
289 populations, and the proportion of inter-population transmissions made up of animal-to-
290 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
291 (blue: A-D) and DTA (red: E-F) models versus the number of isolates sampled from the same
292 outbreak. The circles represent the mean, the error bars represent the 95% HPD interval, the
293 black horizontal lines represent the known parameters for the sampled outbreaks, and the
294 grey horizontal lines represent the known parameters for the entire outbreak.

295

296 **Disproportionate sampling**

297 The DTA and SC models responded to variation in sample proportions for the same
298 simulated outbreak differently. The DTA model's mean estimates showed a much stronger
299 positive correlation with the proportion of isolates sampled from each host population than
300 the SC models' mean estimates (Fig 4). The DTA model's mean estimates displayed a
301 sigmoid-like association with the proportion of isolates sampled from each host population
302 (Fig 5).

303

304 **Fig 4.** Bar graph of the correlation coefficients between the models' mean estimates
305 and the proportion of sampled isolates that are animal or human hosts for the same outbreak
306 that was disproportionately sampled.

307

308 **Fig 5.** The proportion of time spent in the animal (A and E) and human (B and F) host
309 populations, and the proportion of inter-population transmissions made up of animal-to-
310 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
311 (blue: A-D) and DTA (red: E-F) models versus the proportion of sampled isolates that are
312 animal (A, C, E and G) and human (B, D, F and H) for the same outbreak that was
313 disproportionately sampled. The diagonal line represents accurate parameter estimates of the
314 sampled outbreaks, the dots represent the mean, and the error bars represent the 95% HPD
315 interval.

316

317 **Multiple variable simulations**

318 The DTA and SC models showed different associations between known and estimated
319 parameters when 100 isolates were randomly sampled from each of the 23 simulated
320 outbreaks. The SC model predicted a larger proportion of known population and transmission
321 parameters within its 95% HPD interval compared to the DTA model (Fig 6). However, its

322 mean 95% HPD interval sizes were larger and the DTA model's mean estimates showed a
323 stronger positive correlation with the known parameter values than the SC model's mean
324 estimates. Both models had similar mean squared errors between the known parameters and
325 the models' mean estimates. However, the SC model's mean population estimates were all
326 within the 0.2-0.8 interval and its mean transmission rates were all within the 0.35-0.65
327 interval, whilst the DTA models had mean estimates that lay outside of these ranges (Fig 7).

328

329 **Fig 6.** The proportion of outbreak simulations that the models included the known
330 parameter within their 95% highest posterior density (HPD) intervals (A); the correlation
331 coefficient between known parameters and the models' mean estimates (B); the mean squared
332 error between known parameters and the models' mean estimates (C); and the size of the
333 models' 95% HPD intervals (D), for the population and transmission estimates made by the
334 DTA (red) and SC (blue) models for 23 randomly-sampled simulated outbreaks that 100
335 isolates were randomly sampled from.

336

337 **Fig 7.** The proportion of time spent in the animal (A and E) and human (B and F) host
338 populations, and the proportion of inter-population transmissions made up of animal-to-
339 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
340 (blue: A-D) and DTA (red: E-F) models versus the true parameters for 23 simulated
341 outbreaks that 100 isolates were randomly sampled from. The diagonal line represents
342 accurate parameter estimates of the sampled outbreaks, the dots represent the mean, and the
343 error bars represent the 95% HPD interval.

344

345 The phylogenetic trees produced by the DTA and SC models for the 23 simulated
346 outbreaks poorly reflected the sampled transmission trees (Fig 8). The DTA model was

347 unable to detect transmissions along branches in the transmission trees. The SC model could
348 identify transmissions along branches, but often over-estimated the amount of transmissions
349 compared to the true transmission tree. In the example given, the SC model predicted that
350 ‘*Salmonella*’ was predominantly in the animal (red) population, as indicated by the
351 predominantly red branches, but that coalescent events primarily occurred in the human
352 (blue) population. This was common for most of the *maximum a priori* trees produced by the
353 SC model, where the population that was estimated to have a smaller effective population
354 size would be where the coalescent events took place, whilst the population with the
355 estimated larger effective population size would predominate the branches. The phylogenetic
356 trees in Fig 8 represent the most likely trees estimated using the DTA and SC models for one
357 simulated outbreak, not the variation amongst each model, as each model estimated
358 thousands of phylogenetic trees.

359

360 **Fig 8.** Sampled transmission tree (A), maximum clade credibility tree produced by the
361 DTA model (B) and *maximum a posteriori* tree produced by the SC model (C), for one of the
362 23 simulated outbreaks that 100 isolates were randomly sampled from. The blue areas
363 represent time spent in the human population and the red areas represent time spent in the
364 animal population.

365

366 **Equal-time sampling**

367 The DTA and SC models gave similar population and transmission estimates for the
368 23 simulated outbreaks with random (Fig 6-7) and equal-time sampling (Fig 9-10) of 100
369 isolates. Random sampling estimated more known parameters within its 95% HPD interval,
370 but equal-time sampling had smaller mean squared errors between known parameters and the
371 mean estimates, and smaller 95% HPD intervals. The SC and DTA models also estimated

372 similar phylogenetic trees for simulated outbreaks that were sampled using random and
373 equal-time sampling (Fig 11). This suggests that neither sampling method was more suitable
374 for these ancestral state reconstruction models.

375

376 **Fig 9.** The proportion of outbreak simulations that the models included the known
377 parameter within their 95% highest posterior density (HPD) intervals (A); the correlation
378 coefficient between known parameters and the models' mean estimates (B); the mean squared
379 errors between known parameters and the models' mean estimates (C), and the size of the
380 models' 95% HPD intervals (D), for the population and transmission estimates made by the
381 DTA (red) and SC (blue) models for 23 simulated outbreaks that 100 isolates were sampled
382 equally over time from.

383

384 **Fig 10.** The proportion of time spent in the animal (A and E) and human (B and F)
385 host populations, and the proportion of inter-population transmissions made up of animal-to-
386 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
387 (blue: A-D) and DTA (red: E-F) models versus the true parameters for 23 simulated
388 outbreaks that 100 isolates were sampled equally over time from. The diagonal line
389 represents accurate estimates of the sampled outbreaks, the dots represent the mean, and the
390 error bars represent the 95% HPD interval.

391

392 **Fig 11.** Sampled transmission tree (A and D), maximum clade credibility tree
393 produced by the DTA model (B and E) and *maximum a posteriori* tree produced by the SC
394 model (C and F), for one of the 23 simulated outbreaks that 100 isolates were sampled
395 randomly (A-C) and equally over time (D-F). The blue areas represent time spent in the
396 human population and the red areas represent time spent in the animal population.

397

398 **Equal intra-population transmission rates and infectious periods**

399 The DTA and SC models provided more accurate estimates of population parameters
400 for the 12 simulated outbreaks with equal intra-population transmission rates and infectious
401 periods (EPTI) (Fig 12 and 13) than the 23 simulations where these parameters varied (Fig 6
402 and 7), with smaller mean squared errors, a higher proportion of known parameter within
403 their 95% HPD intervals, and mean estimates that were more positively correlated with the
404 known parameters. The DTA model's mean population estimates displayed a sigmoid shape,
405 similar to the simulated outbreak that was disproportionately sampled (Fig 5). On the other
406 hand, the DTA and SC models gave less accurate transmission estimates for the 12 outbreaks
407 with equal intra-population transmission rates and infectious periods between host
408 populations than for the 23 simulations where these parameters varied, with larger mean
409 squared errors, a lower proportion of known parameter within their 95% HPD intervals, and
410 mean estimates that were less positively correlated or negative correlated with the known
411 parameters.

412

413 **Fig 12.** The proportion of outbreak simulations that the models included the known
414 parameter within their 95% highest posterior density (HPD) intervals (A); the correlation
415 coefficients between known parameters and the models' mean estimates (B); the mean
416 squared error between known parameters and the models' mean estimates (C); and the size of
417 the models' 95% HPD intervals (D), for the population and transmission estimates made by
418 the DTA (red) and SC (blue) models for 12 EPTI simulated outbreaks that 100 isolates were
419 randomly sampled from.

420

421 **Fig 13.** The proportion of time spent in the animal (A and E) and human (B and F)
422 host populations, and the proportion of inter-population transmissions made up of animal-to-
423 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
424 (blue: A-D) and DTA (red: E-F) models versus the true parameters for 12 EPTI simulated
425 outbreaks that 100 isolates were randomly sampled from. The diagonal line represents
426 accurate estimates of the sampled outbreaks, the dots represent the mean, and the error bars
427 represent the 95% HPD interval.

428

429 The phylogenetic trees estimated for the 12 EPTI outbreaks (Fig 14) were like those
430 of previous simulated outbreaks (Fig 8). They also demonstrated that the DTA model was
431 unable to estimate ancestral branch states that were a different host population to daughter
432 branches and tips. The SC model could estimate the state of ancestral branches that differed
433 to the tips, but often estimated these branches inaccurately.

434

435 **Fig 14.** Sampled transmission tree (A), maximum clade credibility tree produced by
436 the DTA model (B) and maximum a posteriori tree produced by the SC model (C), for a EPTI
437 simulated outbreak that 100 isolates were randomly sampled from. The blue areas represent
438 time spent in the human population and the red areas represent time spent in the animal
439 population.

440

441 **Host sampling effect on the models' estimates**

442 To determine the effect of host sampling on the SC and DTA models' estimates, the
443 correlation coefficient between the proportion of samples isolated from each host population
444 and the mean estimates for the simulated outbreaks were calculated (Fig 15; S1-S3 Fig). The
445 DTA model's mean population and transmission estimates were more positively correlated

446 with the proportion of samples isolated from each population, than the SC model's. The DTA
447 model's mean estimates displayed similar correlation coefficients for the 12 EPTI simulations
448 and the 23 simulated outbreaks that were sampled randomly and equally over time, whilst the
449 SC model's estimates gave different correlation coefficients for these datasets.

450

451 **Fig 15.** Bar graph of the correlation coefficients between the SC and DTA models'
452 mean estimates and the proportion of isolates sampled from each host population for 12 EPTI
453 simulated outbreaks that 100 isolates were randomly sampled from, and 23 simulated
454 outbreaks that 100 isolates were sampled randomly and equally over time.

455

456 To determine if the difference in sampling fraction could account for the DTA
457 model's estimates for the simulated outbreaks, the correlation coefficient between the
458 proportion of samples isolated from each host and the known parameters were calculated (Fig
459 16; S4-S6 Fig). The known population parameters for the 12 EPTI simulated outbreaks and
460 the sampling proportions were highly correlated, accounting for the more accurate estimates
461 of these known parameters by the DTA model (Fig 13) compared to the known transmission
462 parameters and other outbreak datasets where there was less correlation (Fig 7, 10, 13).

463

464 **Fig 16.** Bar graph of the correlation coefficients between the proportion of isolates
465 sampled from each host population and the known population and transmission parameters
466 for 12 EPTI simulated outbreaks that 100 isolates were randomly sampled from, and 23
467 simulated outbreaks that 100 isolates were sampled randomly and equally over time.

468

469 **DT160 outbreak**

470 The SC and DTA models both predicted that DT160 spent most of the time in the
471 animal host population over the course of the DT160 outbreak in New Zealand (Fig 17).
472 However, the SC model predicted that there were relatively equal amounts of transmission
473 between the animal and human host populations, whilst the DTA model predicted that there
474 was a large amount of animal-to-human transmission and relatively less human-to-animal
475 transmission. The phylogenetic trees estimated for the DT160 outbreak also displayed larger
476 intervals between coalescent events later in the outbreak compared to the outbreaks simulated
477 in this study (Fig 18).

478

479 **Fig 17.** Estimates of the proportion of time spend in the animal (A) and human (B)
480 host populations, and the proportion of inter-population transmissions made up of animal-to-
481 human (C) and human-to-animal (D) transmissions for the DT160 outbreak, as estimated by
482 the SC (blue) and DTA (red) models on 109 isolates. The circles represent the mean and the
483 error bars represent the 95% HPD interval.

484

485 **Fig 18.** Maximum clade credibility tree produced by the DTA model (A) and
486 *maximum a posteriori* tree produced by the SC model (B), based on 109 DT160 isolates.

487

488 **Discussion**

489 The DTA and SC models are ancestral state reconstruction models that were designed
490 to estimate the ancestral state of a group of organisms based on their evolutionary history
491 (3,4). In this study we demonstrated using simulated outbreaks and a previously described
492 salmonellosis outbreak that neither of these models could accurately estimate known
493 population and transmission parameters for these outbreaks.

494 The DTA model assumes that the proportion of samples from each host population is
495 proportional to its relative size (6). This is a problem for outbreaks involving multiple host
496 populations, as the host populations may be sampled at different rates, resulting in samples
497 disproportional to the number of individuals infected within each host population. The
498 simulated outbreaks in this study were stratified by host population before random sampling
499 in efforts to meet this assumption. However, differing intra-population transmission rates and
500 infectious periods between the host populations resulted in inter-population transmission rates
501 and length of times spent in host populations disproportionate to the number of individuals
502 infected within each host population and thus the proportion of each population sampled.
503 This may explain why the DTA model consistently over-estimated the length of time in the
504 animal host population and the number of animal-to-human transmissions for the initial 23
505 simulated outbreaks, as the human host populations of these outbreaks were simulated to
506 have longer infectious periods than the animal host populations. This resulted in longer
507 periods spent in the human host population and a larger number of human-to-animal
508 transmissions relative to the number of humans sampled.

509 The DTA model appeared to estimate population parameters more accurately when
510 the parameter was directly proportional to the number of isolates from each host population
511 sampled. In these instances, the population estimates and simulated outbreak parameters
512 shared a sigmoid-like relationship due to the model's ancestral branch estimates: the DTA
513 model usually predicts that all the ancestral branches are one host population, until the
514 majority of the tips are another host population, where all the ancestral branches switch (11).
515 The correct population parameters were also only estimated when simulating outbreaks with
516 equal intra-population transmission rates and infectious periods, parameters that usually
517 differ between *Salmonella* host populations (15,16). However, even in these instances the

518 DTA model inaccurately estimated ancestral host population states and transmission
519 parameters.

520 The SC model gave similar estimates for all the simulated outbreaks. It was poor at
521 estimating simulated outbreaks known parameters, only accurately estimating them when
522 they were within the range that it consistently estimated. The SC model's inaccurate
523 estimates are possibly due to the model's assumption that the effective population size of the
524 host populations were consistent throughout the outbreak (10), which does not apply to
525 salmonellosis outbreaks whose effective population size varies over the course of the
526 outbreak (11). There may be other reasons why the SC model was unable to detect a signal,
527 but it is difficult to test for these without first accounting for the model's effective population
528 size assumption.

529 The inability of the SC and DTA models to accurately estimate salmonellosis
530 outbreak parameters highlights the need for outbreak-specific models. These models would
531 need to be able to take into consideration variable sampling between host populations, like
532 the SC model, and changes in the effective population size, like the DTA model. In addition,
533 they would need to be able to take into consideration variation in infectious periods and intra-
534 population transmission rates.

535 The MASTER package of BEAST2 allowed many salmonellosis outbreaks to be
536 simulated using the stochastic SIR model. The simulated outbreaks contained a large amount
537 of variation in the amount of time spent in the animal and human host populations, but less
538 variation in inter-population transmissions due to only simulating two host populations.
539 Therefore, unequal transmission values were only simulated using one very high and one
540 very low inter-population transmission value. This in part explains why the SC model was
541 more likely to provide estimates that matched known simulation parameters because it always
542 gave similar mean estimates around the 0.35-0.65 range, which most of the known

543 transmission parameters for the simulated outbreaks were within. Further work with multiple
544 host populations may help better understand these models' application to salmonellosis
545 outbreaks.

546 The DTA and SC models' estimates of the DT160 outbreak underline some of the
547 limitations of this study. The DTA model estimated that DT160 spent most of its time in the
548 animal host population and that there was a larger amount of animal-to-human transmission
549 than human-to-animal transmission, which is to be expected as the DTA model is affected by
550 sample size and a larger number of animal isolates were analyzed than human isolates in the
551 DT160 study. The SC model estimated similar amounts of animal-to-human transmission
552 than human-to-animal transmission, which is also to be expected as our study shows it
553 usually gives similar transmission rates between two host populations. However, the SC
554 model estimated that DT160 spent over 90% of its time in the animal host population and less
555 than 10% of its time in the human host population, outside the 20-80% range estimated for
556 simulated outbreaks, and both models produced phylogenetic trees with larger distances
557 between coalescent events towards the later part of the outbreak than simulated outbreaks.
558 The effective population size affects the timing of coalescent events for randomly sampled
559 individuals (17). This suggests that the DT160 outbreak had a much larger effective
560 population size than any of the simulated outbreaks in this study. It also indicates that the SC
561 model's estimates maybe influenced by branch length. Simulations with larger effective
562 population sizes are required to test this.

563 In conclusion, our comparison of applicability of the SC and DTA models to
564 salmonellosis outbreaks between the known parameters of simulated outbreaks and the
565 models' estimates suggest neither model is appropriate for this analysis. Our findings
566 highlight the need for outbreak-specific models that can also take into consideration intra-

567 population transmission rates, infectious periods, disproportionate sampling and changes in
568 the effective population size.

569

570 **Acknowledgements**

571 We acknowledge the contribution of the New Zealand eScience Infrastructure (NeSI)
572 high-performance computing facilities to the results of this research.

573

574 **References**

- 575 1. WHO. Disease outbreaks [Internet]. World Health Organization. 2016. p. 1–1.
576 Available from: http://www.who.int/topics/disease_outbreaks/en/
- 577 2. Mather AE, Reid SWJ, Maskell DJ, Parkhill J, Fookes MC, Harris S., et al.
578 Distinguishable epidemics of multidrug-resistant *Salmonella* Typhimurium DT104 in
579 different hosts. *Science* (80-) [Internet]. 2013;341(6153):1514–7. Available from:
580 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84885656011&partnerID=40&md5=31f2c7923584046fd043e1f0ebf66e6a)
581 [84885656011&partnerID=40&md5=31f2c7923584046fd043e1f0ebf66e6a](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84885656011&partnerID=40&md5=31f2c7923584046fd043e1f0ebf66e6a)
- 582 3. Lemey P, Rambaut A, Drummond A, Suchard M. Bayesian phylogeography finds its
583 roots. *PLoS Comput Biol* [Internet]. 2009;5(9):1–16. Available from:
584 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-70349675664&partnerID=40&md5=ef46381e24c435e991da9c5654f6bc3b)
585 [70349675664&partnerID=40&md5=ef46381e24c435e991da9c5654f6bc3b](https://www.scopus.com/inward/record.uri?eid=2-s2.0-70349675664&partnerID=40&md5=ef46381e24c435e991da9c5654f6bc3b)
- 586 4. Vaughan TG, Kühnert D, Poppinga A, Welch D, Drummond AJ. Efficient Bayesian
587 inference under the structured coalescent. *Bioinformatics* [Internet].
588 2014;30(16):2272–9. Available from:
589 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84906255301&partnerID=40&md5=3e9fd7aac4d0f3e2f7573b99d56a0136)
590 [84906255301&partnerID=40&md5=3e9fd7aac4d0f3e2f7573b99d56a0136](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84906255301&partnerID=40&md5=3e9fd7aac4d0f3e2f7573b99d56a0136)
- 591 5. Felsenstein J. Evolutionary trees from DNA sequences: A maximum likelihood

- 592 approach. *J Mol Evol* [Internet]. 1981;17(6):368–76. Available from:
593 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-0019797407&doi=10.1007%2FBBF01734359&partnerID=40&md5=eb071314e58843b437de4d67c4e77341)
594 [0019797407&doi=10.1007%2FBBF01734359&partnerID=40&md5=eb071314e58843b](https://www.scopus.com/inward/record.uri?eid=2-s2.0-0019797407&doi=10.1007%2FBBF01734359&partnerID=40&md5=eb071314e58843b437de4d67c4e77341)
595 [437de4d67c4e77341](https://www.scopus.com/inward/record.uri?eid=2-s2.0-0019797407&doi=10.1007%2FBBF01734359&partnerID=40&md5=eb071314e58843b437de4d67c4e77341)
- 596 6. De Maio N, Wu CH, O’Reilly KM, Wilson D. New routes to phylogeography: A
597 Bayesian structured coalescent approximation. *PLoS Genet* [Internet]. 2015;11(8):1–
598 22. Available from: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84940739133&partnerID=40&md5=ac9b1b1462c622ca95b13159ba9b0302)
599 [84940739133&partnerID=40&md5=ac9b1b1462c622ca95b13159ba9b0302](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84940739133&partnerID=40&md5=ac9b1b1462c622ca95b13159ba9b0302)
- 600 7. Gould LH, Walsh KA, Vieira AR, Herman K, Williams IT, Hall AJ, et al. Surveillance
601 for foodborne disease outbreaks - United States, 1998-2008. *MMWR Surveill Summ*
602 [Internet]. 2013;62(1):1–34. Available from:
603 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84893354227&partnerID=40&md5=e446fa7f0c9956ecba67bb3392167787)
604 [84893354227&partnerID=40&md5=e446fa7f0c9956ecba67bb3392167787](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84893354227&partnerID=40&md5=e446fa7f0c9956ecba67bb3392167787)
- 605 8. Wikswo M. Outbreaks of acute gastroenteritis transmitted by person-to-person contact-
606 United States, 2009-2010. *Am J Public Health* [Internet]. 2014;104(11):e13–4.
607 Available from: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84986880401&partnerID=40&md5=d3ff522e8cb9145941b43b13a7253817)
608 [84986880401&partnerID=40&md5=d3ff522e8cb9145941b43b13a7253817](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84986880401&partnerID=40&md5=d3ff522e8cb9145941b43b13a7253817)
- 609 9. Vaughan TG, Drummond AJ. A stochastic simulator of birth-death master equations
610 with application to phylodynamics. *Mol Biol Evol* [Internet]. 2013;30(6):1480–93.
611 Available from: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84877761231&doi=10.1093%2Fmolbev%2Fmst057&partnerID=40&md5=1a2586ee7cf353efea6d6470c3855804)
612 [84877761231&doi=10.1093%2Fmolbev%2Fmst057&partnerID=40&md5=1a2586ee7](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84877761231&doi=10.1093%2Fmolbev%2Fmst057&partnerID=40&md5=1a2586ee7cf353efea6d6470c3855804)
613 [cf353efea6d6470c3855804](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84877761231&doi=10.1093%2Fmolbev%2Fmst057&partnerID=40&md5=1a2586ee7cf353efea6d6470c3855804)
- 614 10. Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, et al. BEAST 2: A
615 Software Platform for Bayesian Evolutionary Analysis. *PLoS Comput Biol* [Internet].
616 2014;10(4). Available from: <https://www.scopus.com/inward/record.uri?eid=2-s2.0->

- 617 84901305512&doi=10.1371%2Fjournal.pcbi.1003537&partnerID=40&md5=5515082
618 bdf0d92199f98e1acae7385fb
- 619 11. Bloomfield SJ, Benschop J, Biggs PJ, Marshall JC, Hayman DTS, Carter PE, et al.
620 Genomic analysis of *Salmonella enterica* serovar Typhimurium DT160 associated
621 with a 14-year outbreak, New Zealand, 1998-2012. *Emerg Infect Dis.* 2017;23(6):906–
622 13.
- 623 12. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with
624 BEAUti and the BEAST 1.7. *Mol Biol Evol* [Internet]. 2012;29(8):1969–73. Available
625 from: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84864530626&partnerID=40&md5=3fe0a37300c038c06df1307c8f1c69d9)
626 [84864530626&partnerID=40&md5=3fe0a37300c038c06df1307c8f1c69d9](https://www.scopus.com/inward/record.uri?eid=2-s2.0-84864530626&partnerID=40&md5=3fe0a37300c038c06df1307c8f1c69d9)
- 627 13. Tavaré S. Some probabilistic and statistical problems in the analysis of DNA
628 sequences. *Am Math Soc.* 1986;17:57–86.
- 629 14. Minin VN., Bloomquist E, Suchard M. Smooth skyride through a rough skyline:
630 Bayesian coalescent-based inference of population dynamics. *Mol Biol Evol* [Internet].
631 2008;25(7):1459–71. Available from:
632 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-45849122039&partnerID=40&md5=83fbb3d858a0cfc4a5ef5aec2df2447)
633 [45849122039&partnerID=40&md5=83fbb3d858a0cfc4a5ef5aec2df2447](https://www.scopus.com/inward/record.uri?eid=2-s2.0-45849122039&partnerID=40&md5=83fbb3d858a0cfc4a5ef5aec2df2447)
- 634 15. Alexander KA, Warnick LD, Cripps CJ, McDonough PL, Grohn YT, Wiedmann M, et
635 al. Fecal shedding of, antimicrobial resistance in, and serologic response to *Salmonella*
636 Typhimurium in dairy calves. *J Am Vet Med Assoc* [Internet]. 2009;235(6):739–48.
637 Available from: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-70349632969&doi=10.2460%2Fjavma.235.6.739&partnerID=40&md5=a99e961243b2435e87df7ce96b1ced42)
638 [70349632969&doi=10.2460%2Fjavma.235.6.739&partnerID=40&md5=a99e961243b](https://www.scopus.com/inward/record.uri?eid=2-s2.0-70349632969&doi=10.2460%2Fjavma.235.6.739&partnerID=40&md5=a99e961243b2435e87df7ce96b1ced42)
639 [2435e87df7ce96b1ced42](https://www.scopus.com/inward/record.uri?eid=2-s2.0-70349632969&doi=10.2460%2Fjavma.235.6.739&partnerID=40&md5=a99e961243b2435e87df7ce96b1ced42)
- 640 16. Murase T, Yamada M, Muto T, Matsushima A, Yamai S. Fecal excretion of
641 *Salmonella enterica* serovar Typhimurium following a food-borne outbreak. *J Clin*

642 Microbiol [Internet]. 2000;38(9):3495–7. Available from:
643 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-0033831284&partnerID=40&md5=d957dabb4032b53db984fc6877b83ccf)
644 [0033831284&partnerID=40&md5=d957dabb4032b53db984fc6877b83ccf](https://www.scopus.com/inward/record.uri?eid=2-s2.0-0033831284&partnerID=40&md5=d957dabb4032b53db984fc6877b83ccf)
645 17. Heled J, Drummond AJ. Bayesian inference of population size history from multiple
646 loci. BMC Evol Biol [Internet]. 2008;8(1):1–15. Available from:
647 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-60549111699&partnerID=40&md5=181e2952d1c48734d15ea5675d2f6327)
648 [60549111699&partnerID=40&md5=181e2952d1c48734d15ea5675d2f6327](https://www.scopus.com/inward/record.uri?eid=2-s2.0-60549111699&partnerID=40&md5=181e2952d1c48734d15ea5675d2f6327)

649

650 **S1 Appendix - Simulated outbreak parameters**

651

652 **S1 Fig.** The proportion of time spent in the animal (A and E) and human (B and F)
653 host populations, and the proportion of inter-population transmissions made up of animal-to-
654 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
655 (blue: A-D) and DTA (red: E-F) models versus the proportion of samples made up of animal
656 (A, C, E and G) and human (B, D, F and H) host populations for 12 EPTI simulated
657 outbreaks that 100 isolates were randomly sampled from. The dots represent the mean, and
658 the error bars represent the 95% HPD interval.

659

660 **S2 Fig.** The proportion of time spent in the animal (A and E) and human (B and F)
661 host populations, and the proportion of inter-population transmissions made up of animal-to-
662 human (C and G) and human-to-animal (D and H) transmissions as estimated by the SC
663 (blue: A-D) and DTA (red: E-F) models versus the proportion of samples made up of animal
664 (A, C, E and G) and human (B, D, F and H) host populations for 23 simulated outbreaks that
665 100 isolates were randomly sampled from. The dots represent the mean, and the error bars
666 represent the 95% HPD interval.

667

668 **S3 Fig.** Scatterplots of the proportion of time spent in the animal (A and E) and
669 human (B and F) host populations, and the proportion of inter-population transmissions made
670 up of animal-to-human (C and G) and human-to-animal (D and H) transmissions as estimated
671 by the SC (blue: A-D) and DTA (red: E-F) models versus the proportion of samples made up
672 of animal (A, C, E and G) and human (B, D, F and H) host populations for 23 simulated
673 outbreaks that 100 isolates were sampled equally over time from. The dots represent the
674 mean, and the error bars represent the 95% HPD interval.

675

676 **S4 Fig.** The proportion of samples made up of animal (A and C) and human (B and
677 D) host populations, versus the known population (A and B) and transmission (C and D)
678 parameters for 12 EPTI simulated outbreaks that 100 isolates were randomly sampled from.

679

680 **S5 Fig.** The proportion of samples made up of animal (A and C) and human (B and
681 D) host populations, versus the known population (A and B) and transmission (C and D)
682 parameters for 23 simulated outbreaks that 100 isolates were randomly sampled from.

683

684 **S6 Fig.** The proportion of samples made up of animal (A and C) and human (B and
685 D) host populations, versus the known population (A and B) and transmission (C and D)
686 parameters for 23 simulated outbreaks that 100 isolates were sampled equally over time from.

Outbreak simulation:
Normal (n=23)
EPTI (n=12)

Transmission tree

Random sampling:
Normal (n=23)
EPTI (n=12)
Model consistency (n=10)
Varying sample size (n=12)

Equal-time sampling:
Normal (n=23)

Disproportionate sampling:
Normal (n=10)

Sampled tree

Sequence simulation
(n=90)

SNPs

DTA model
(n=90)

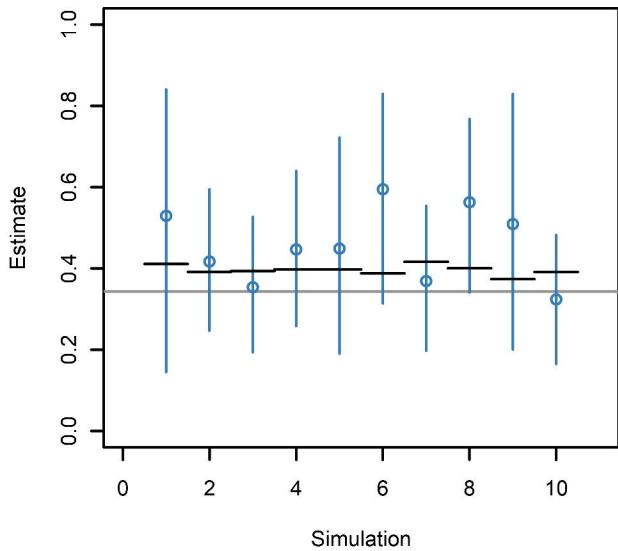
SC model
(n=90)

Log files

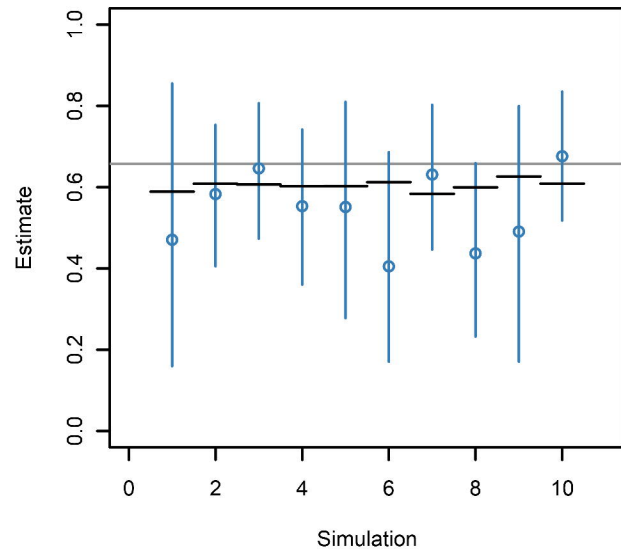
Model comparison
(n=90)

Population

A

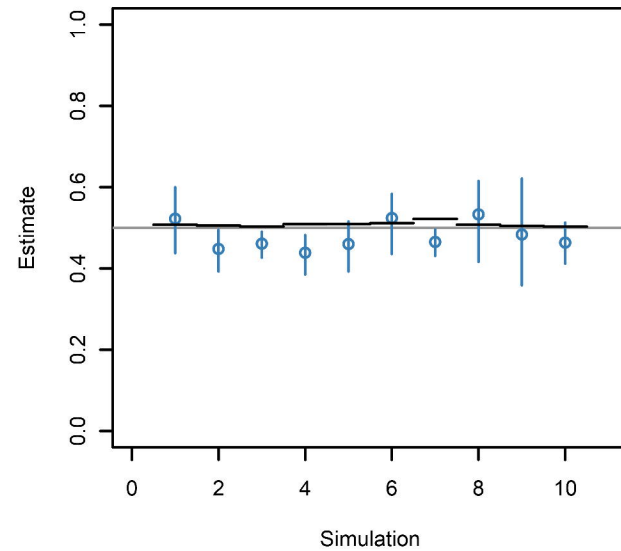


B

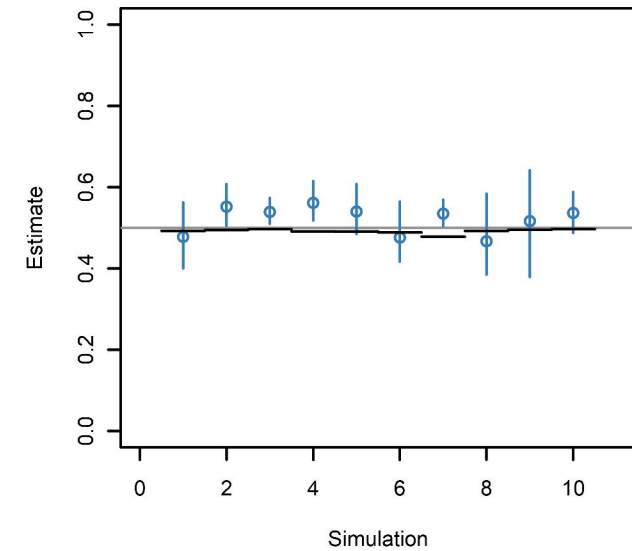


Transmission

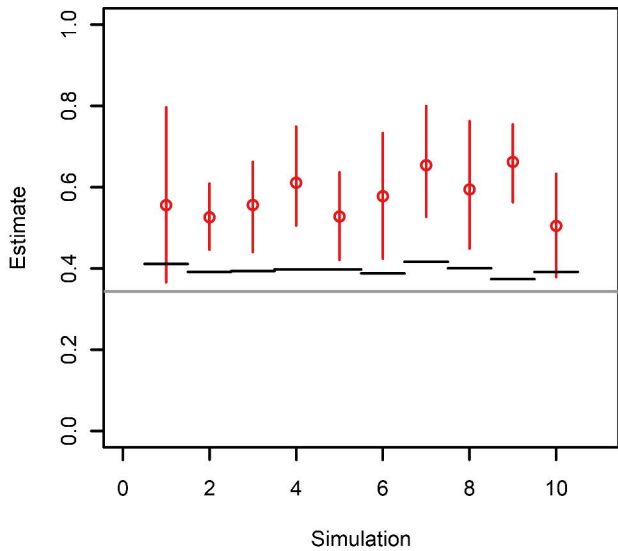
C



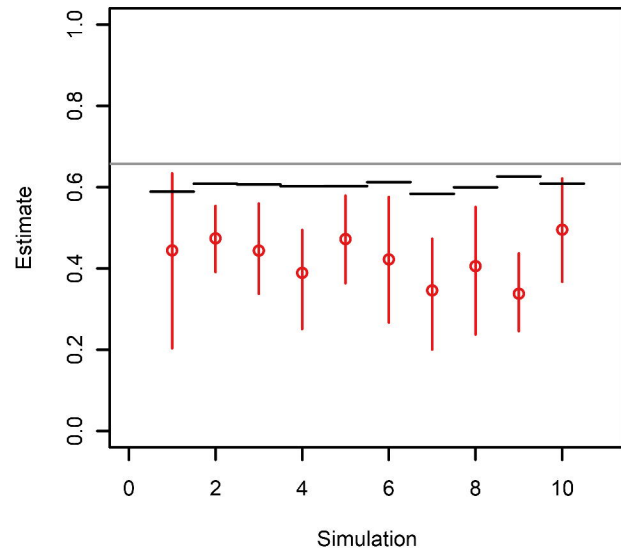
D



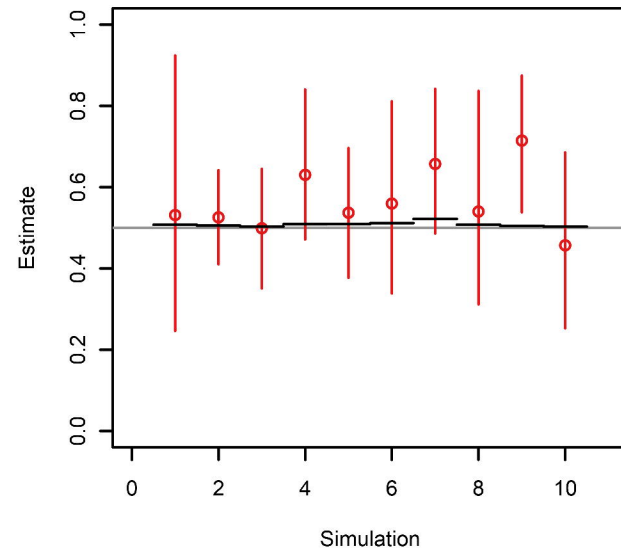
E



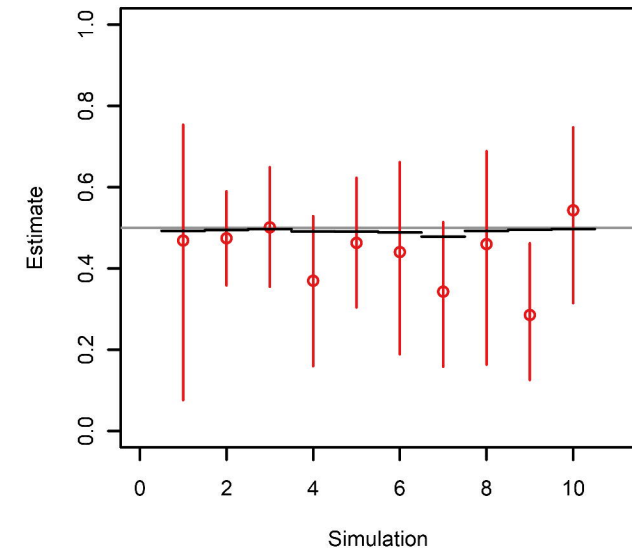
F



G

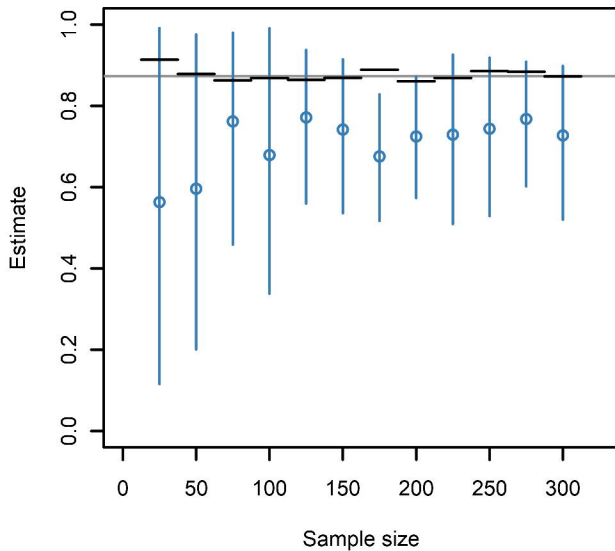


H

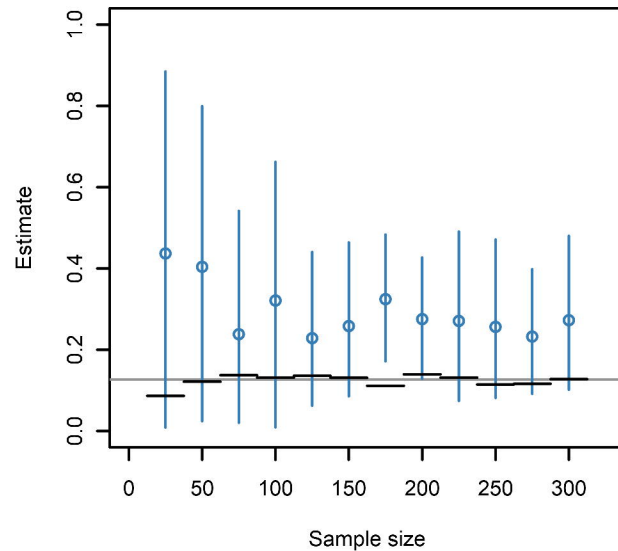


Population

A

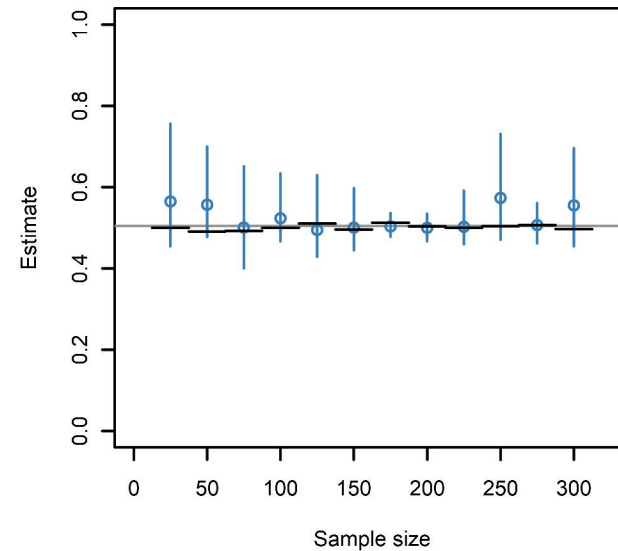


B

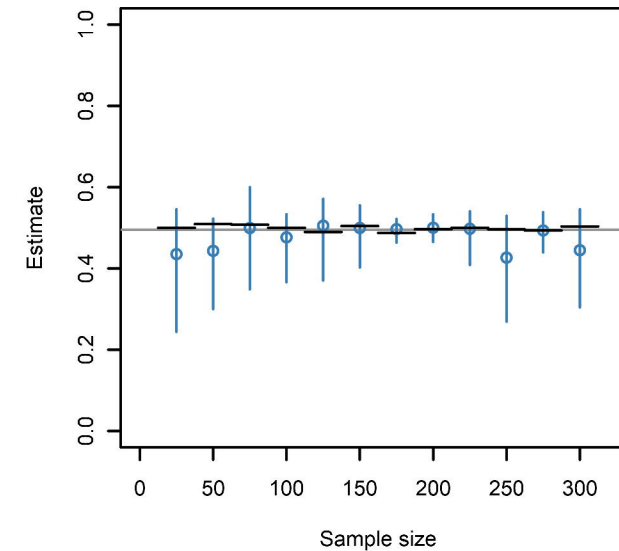


Transmission

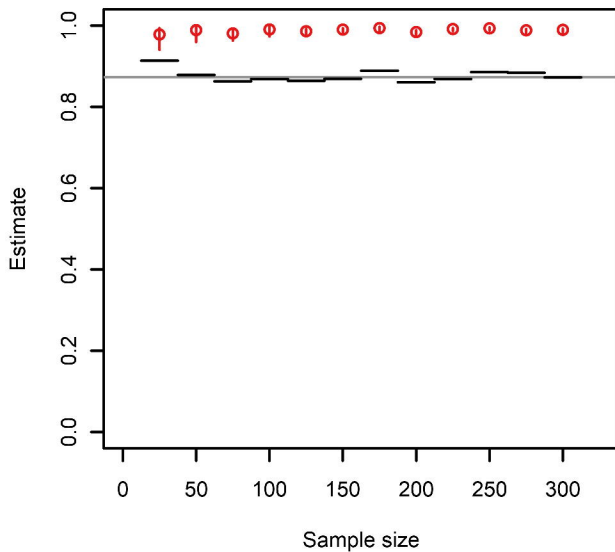
C



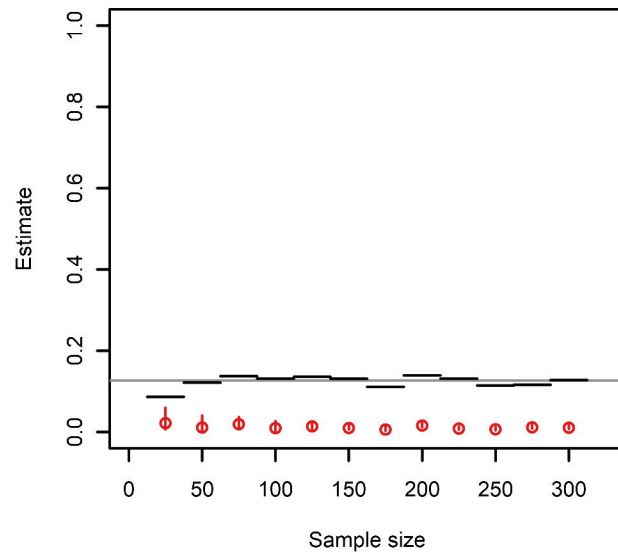
D



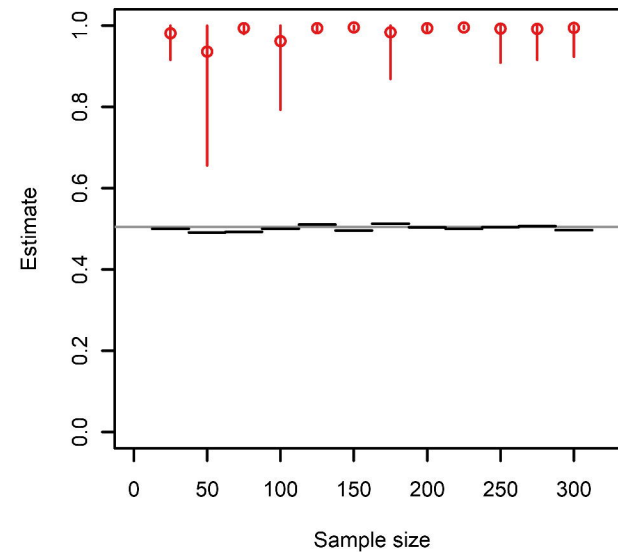
E



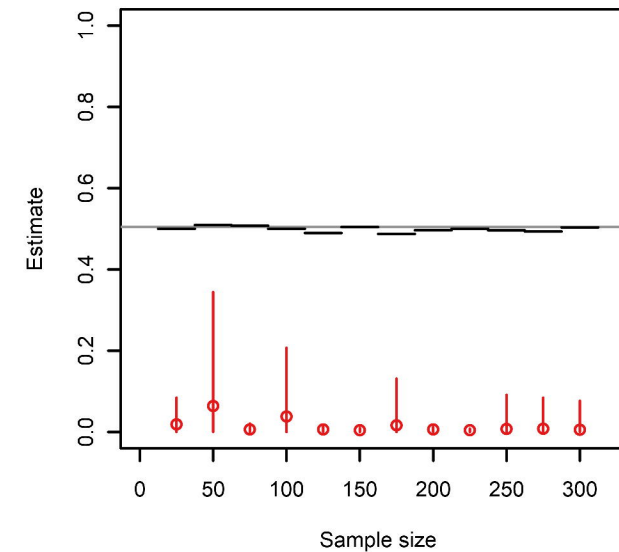
F

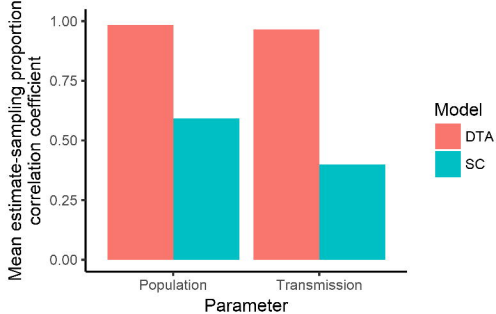


G



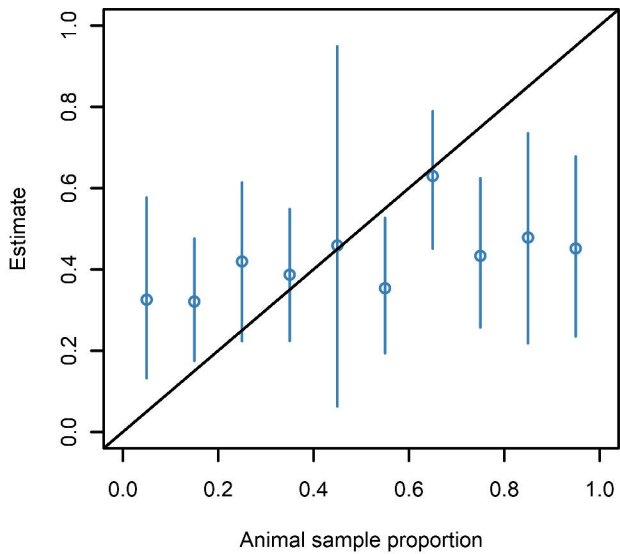
H



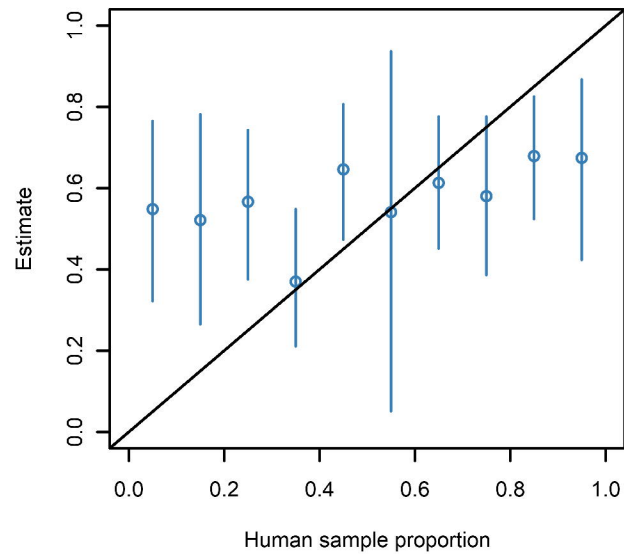


Population

A

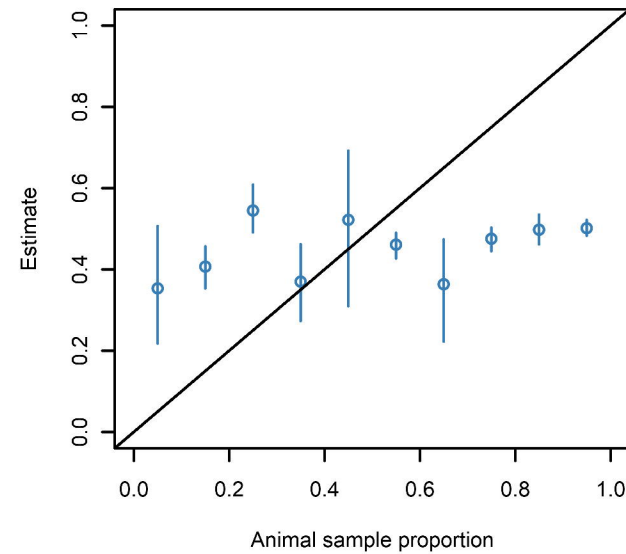


B

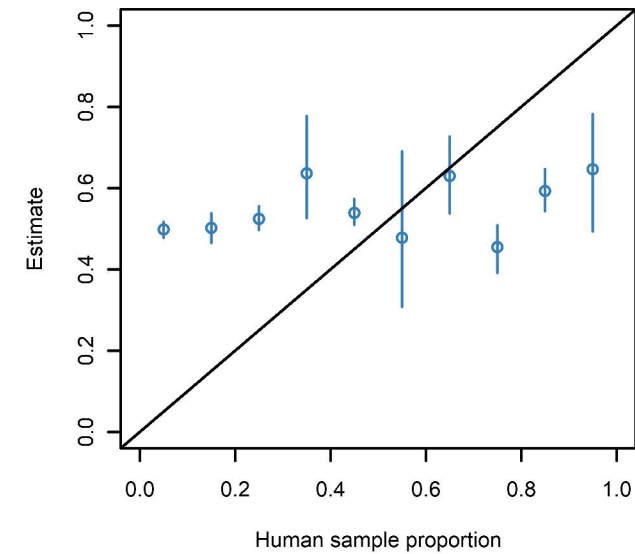


Transmission

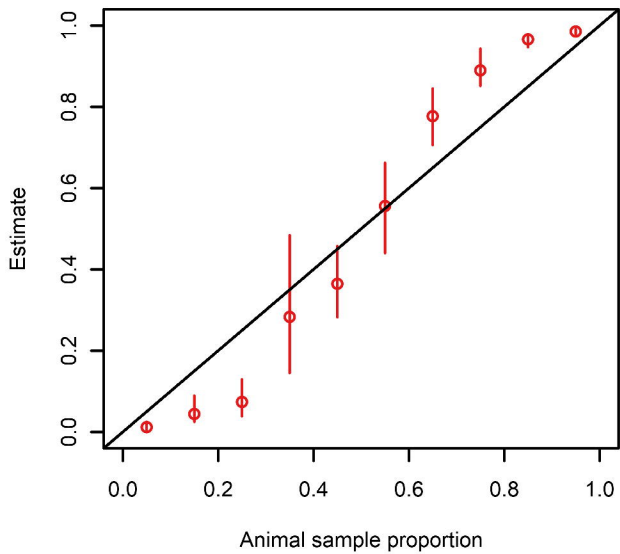
C



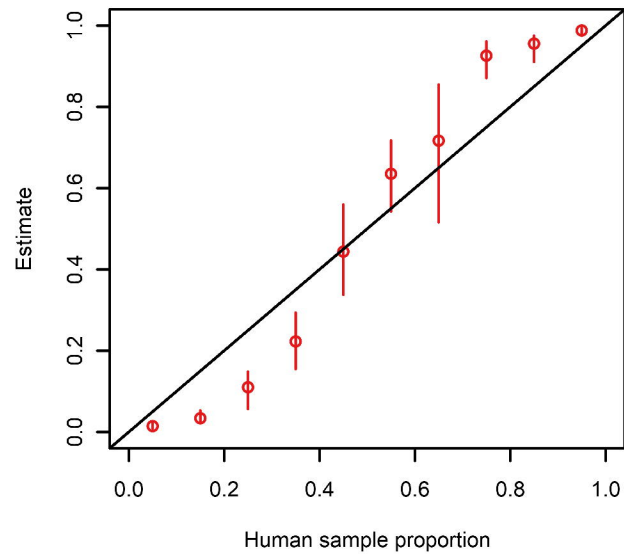
D



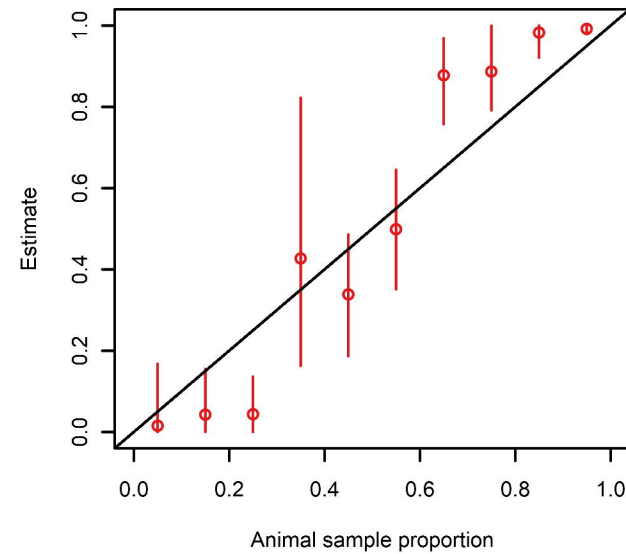
E



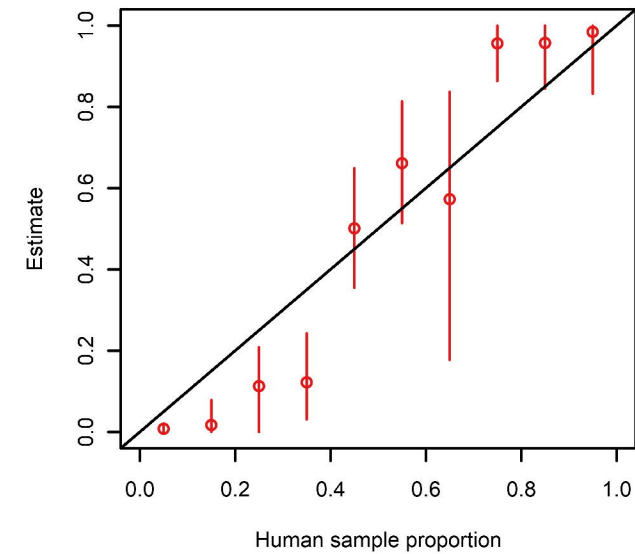
F

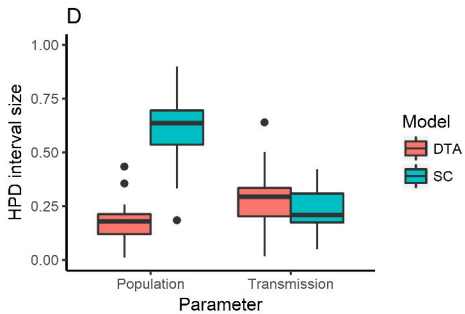
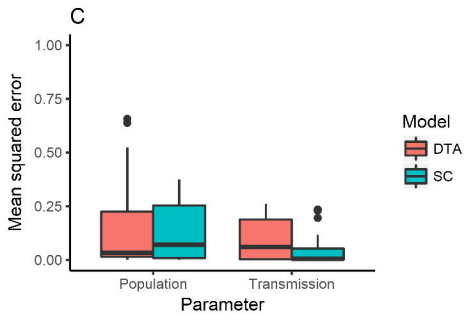
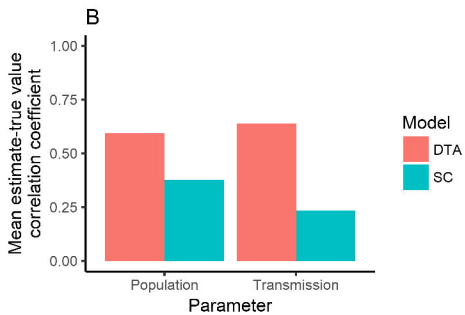
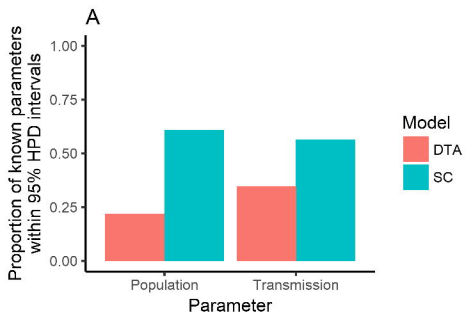


G



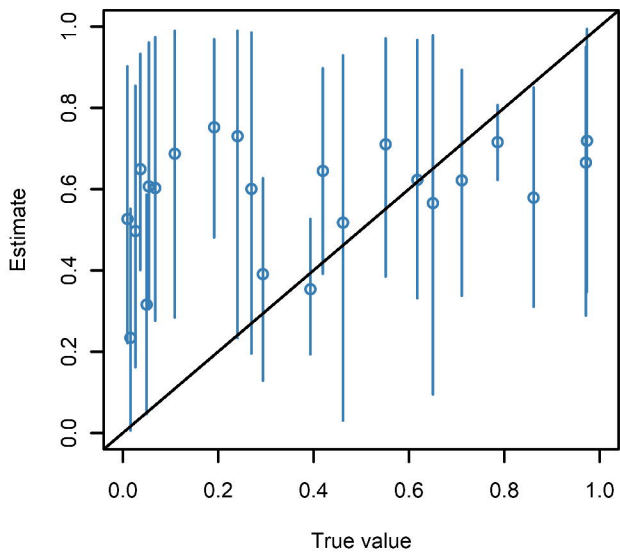
H



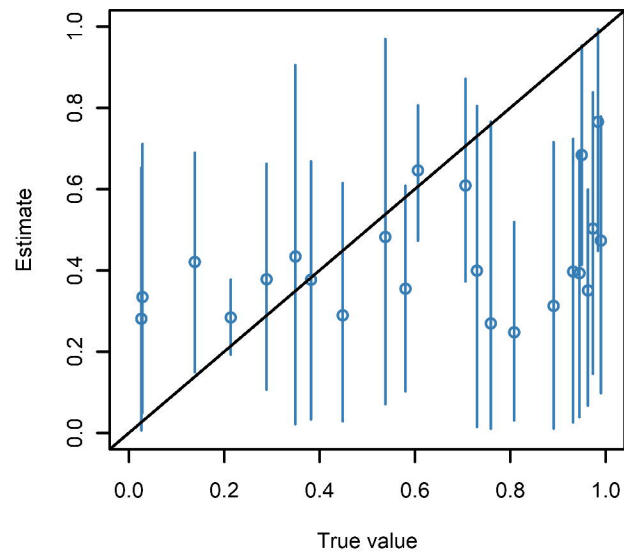


Population

A

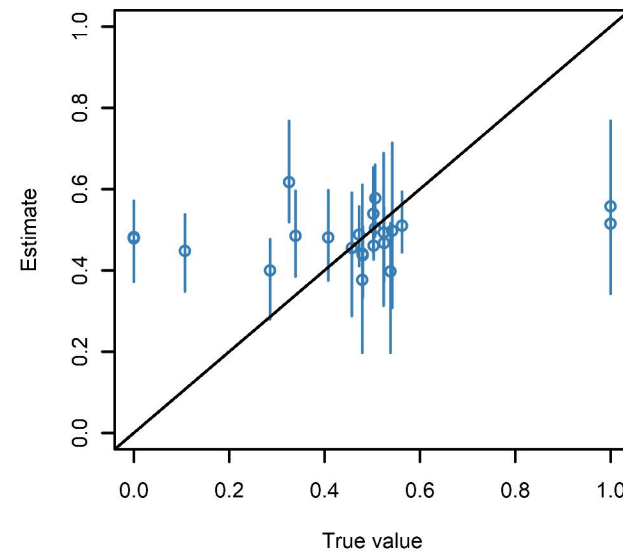


B

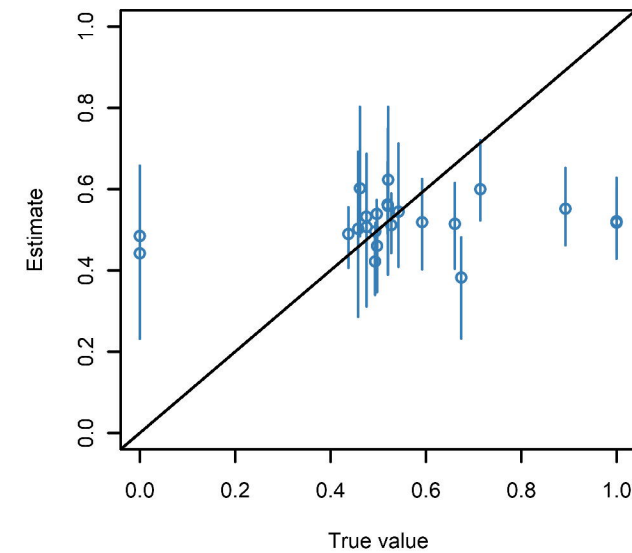


Transmission

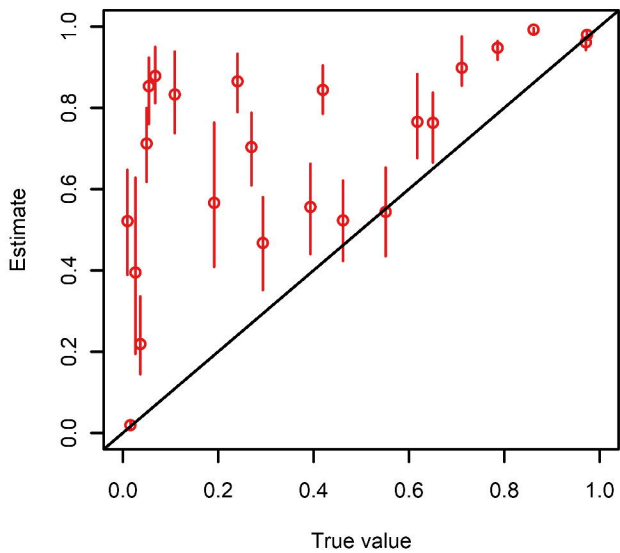
C



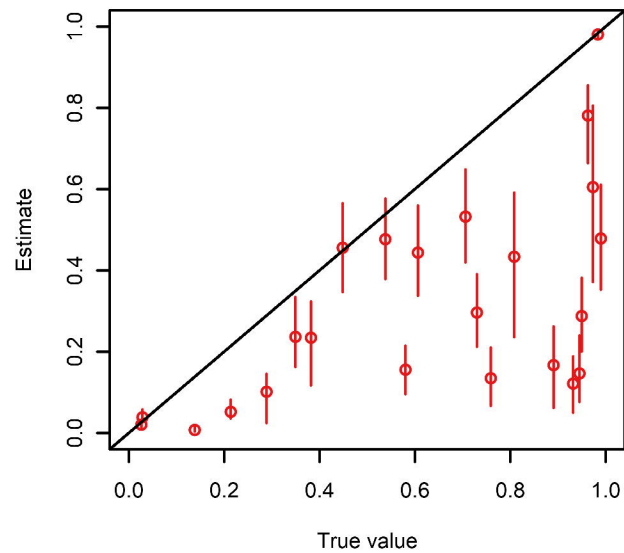
D



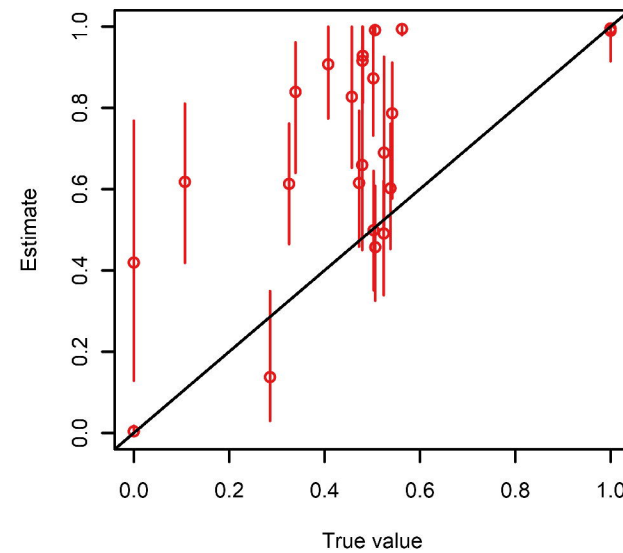
E



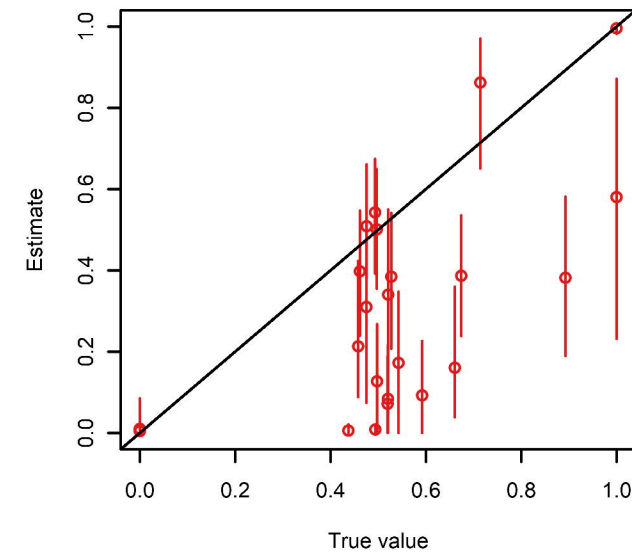
F

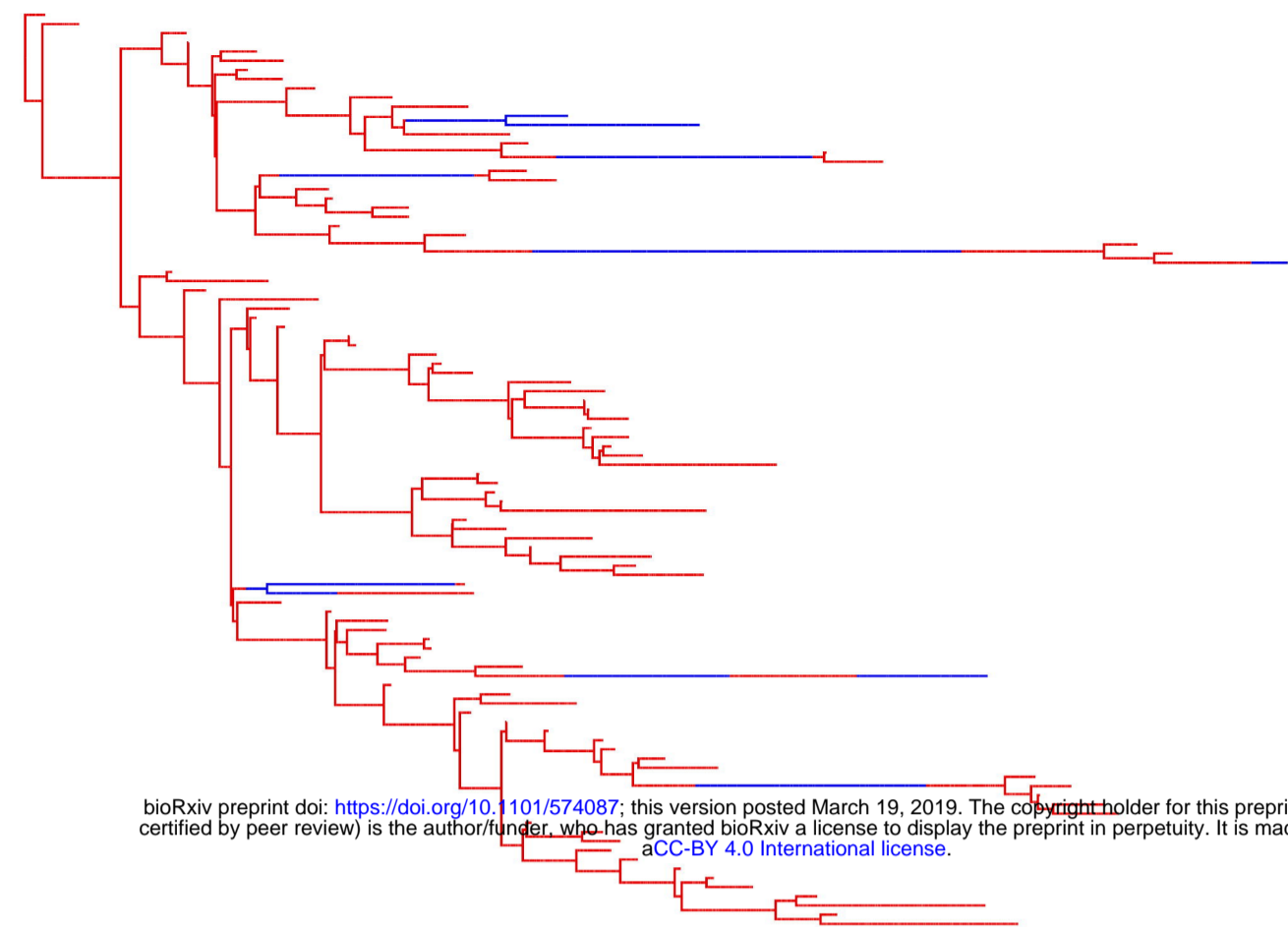


G

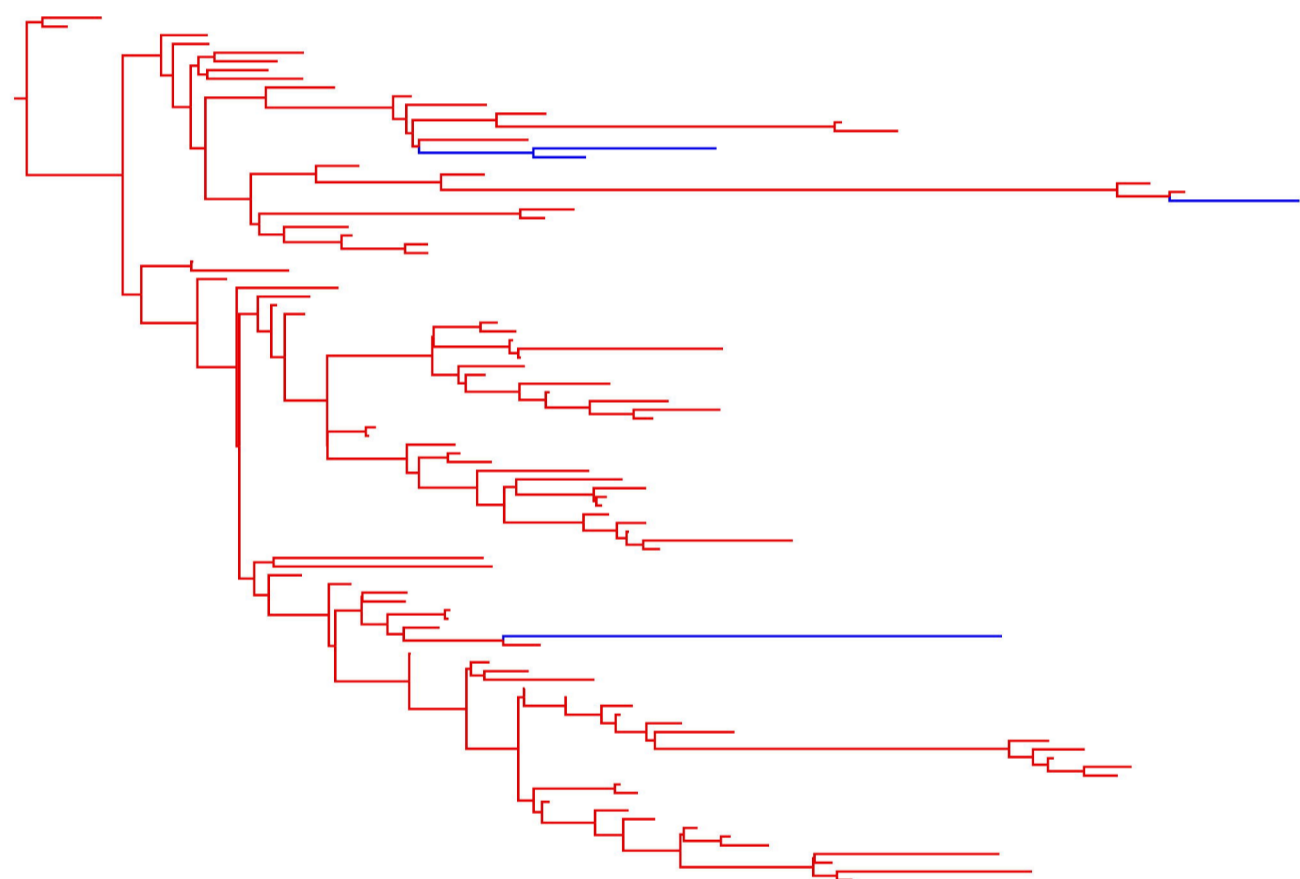
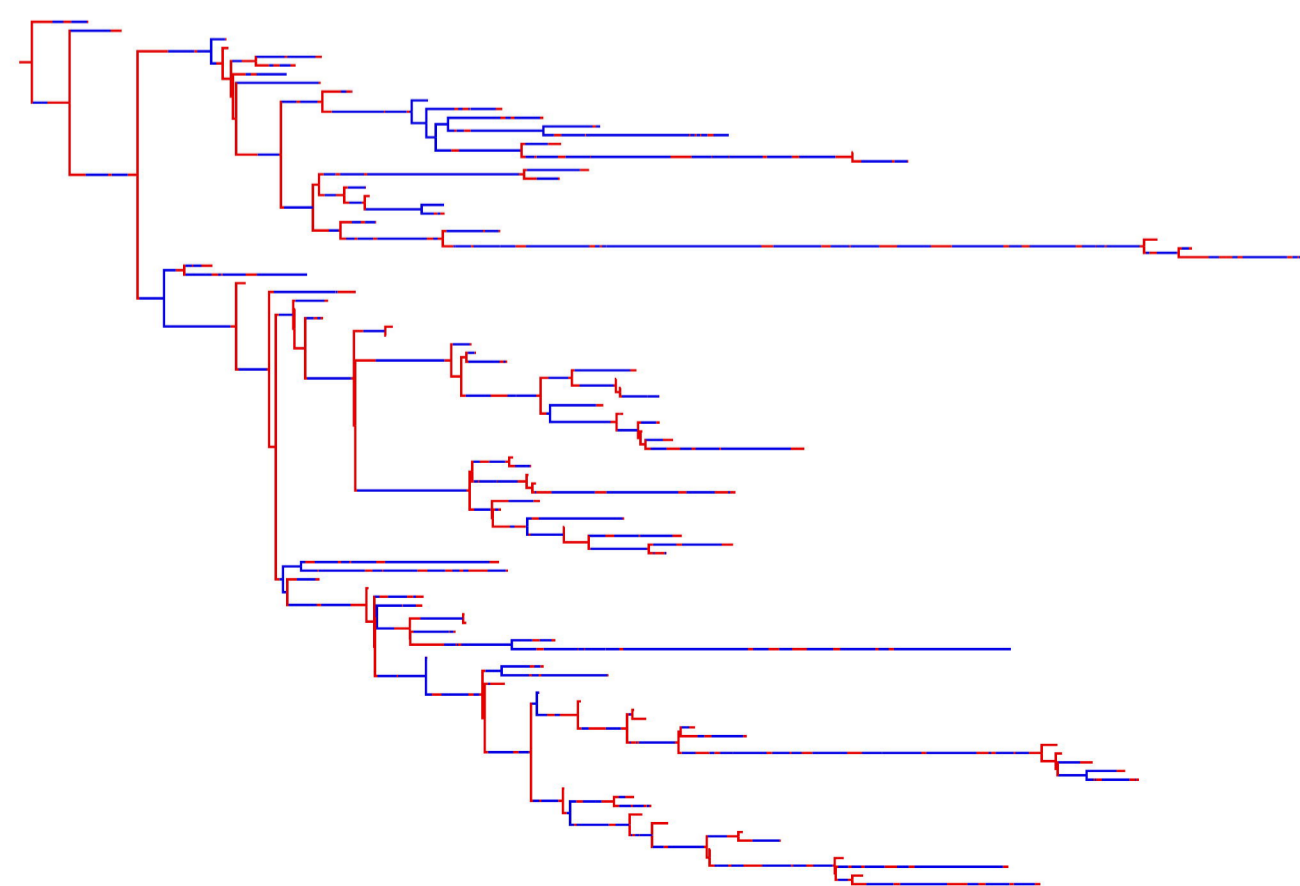


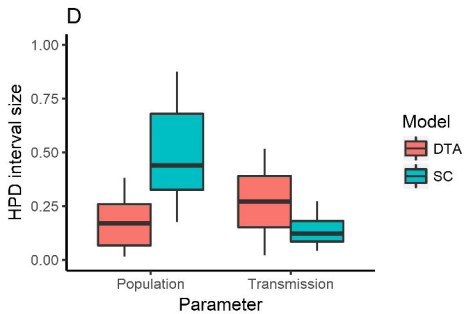
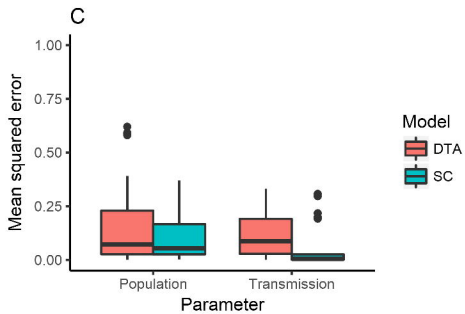
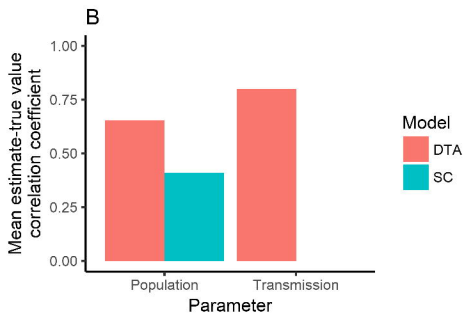
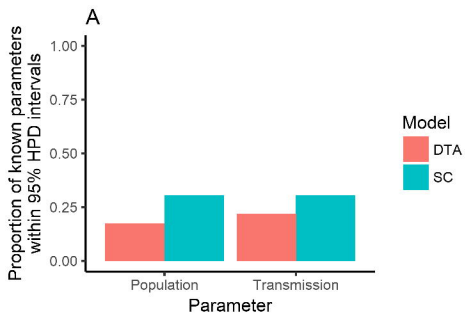
H



A

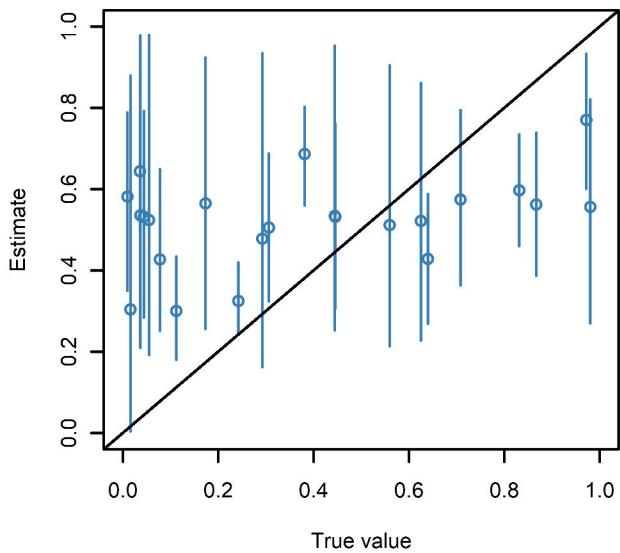
bioRxiv preprint doi: <https://doi.org/10.1101/574087>; this version posted March 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

B**C**

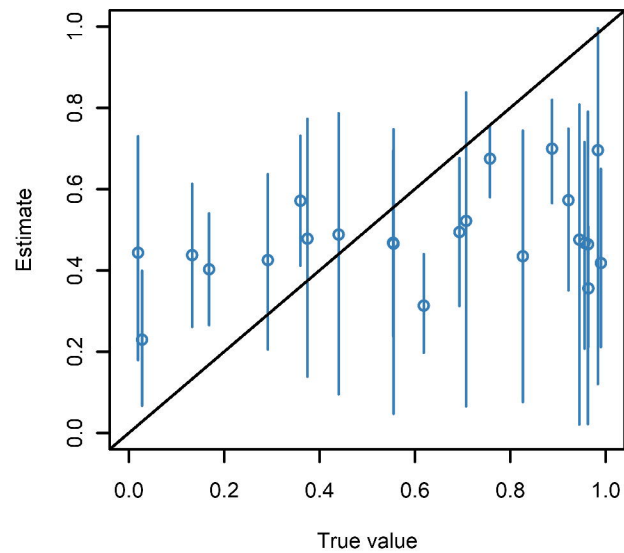


Population

A

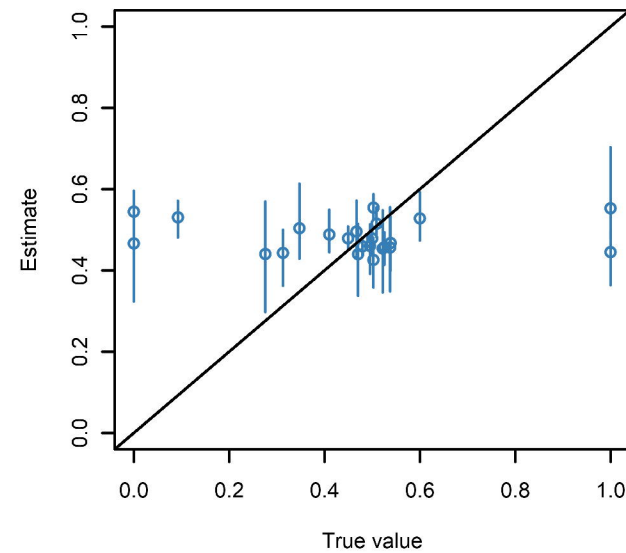


B

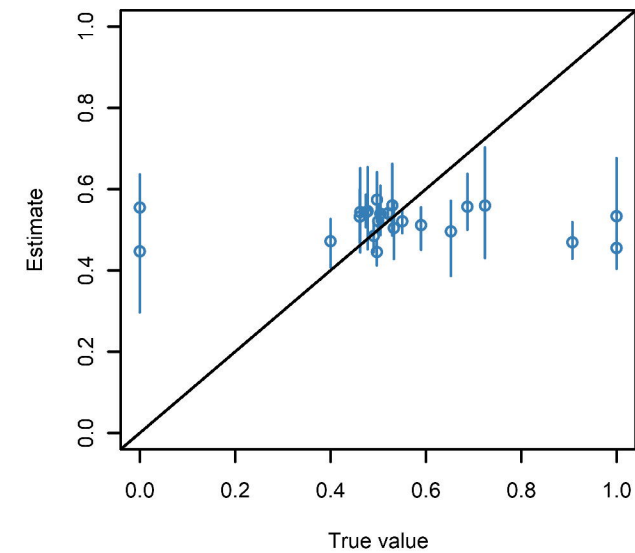


Transmission

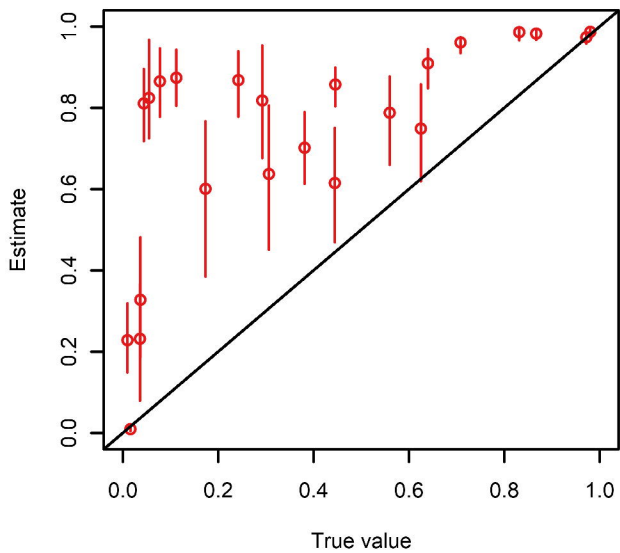
C



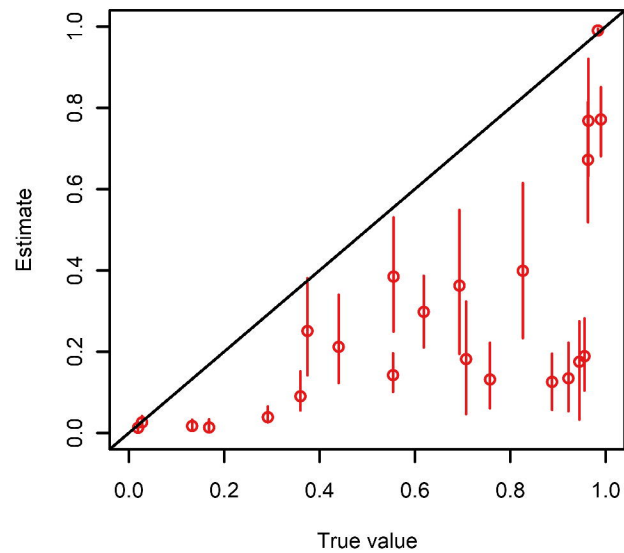
D



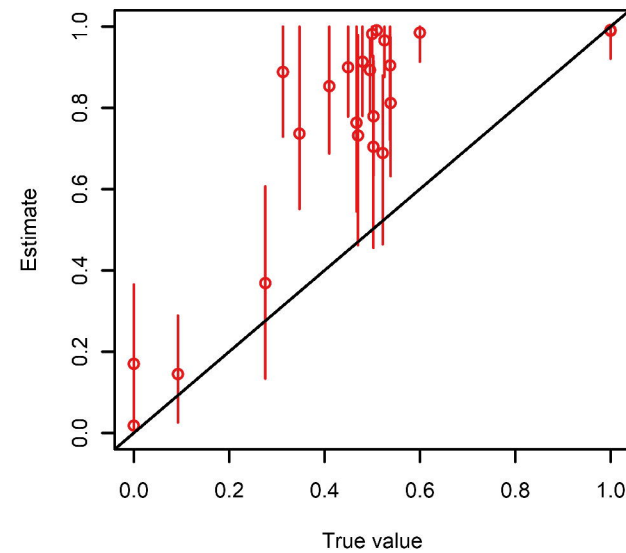
E



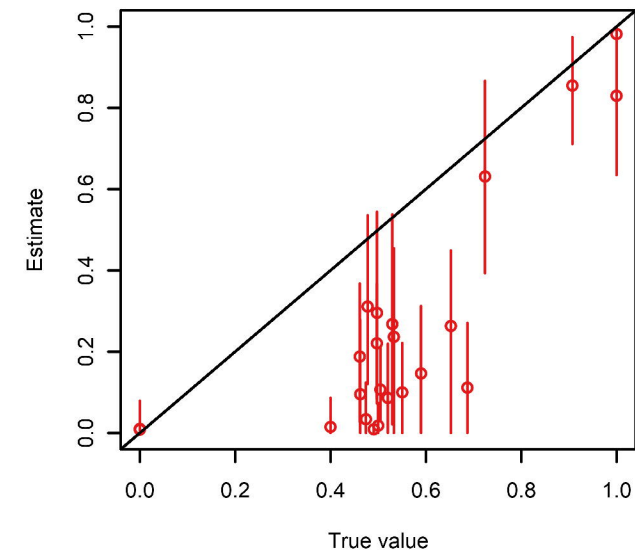
F

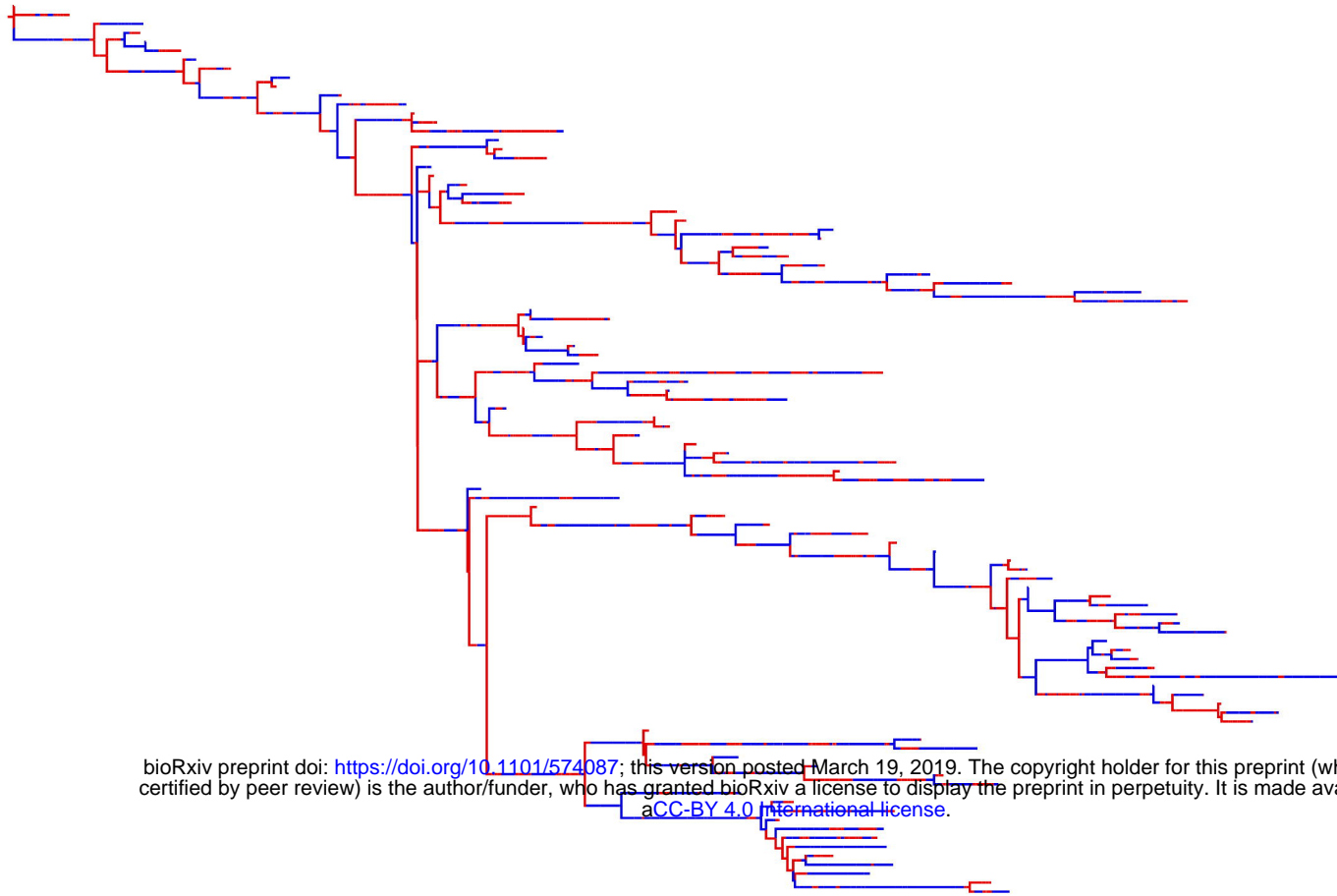


G

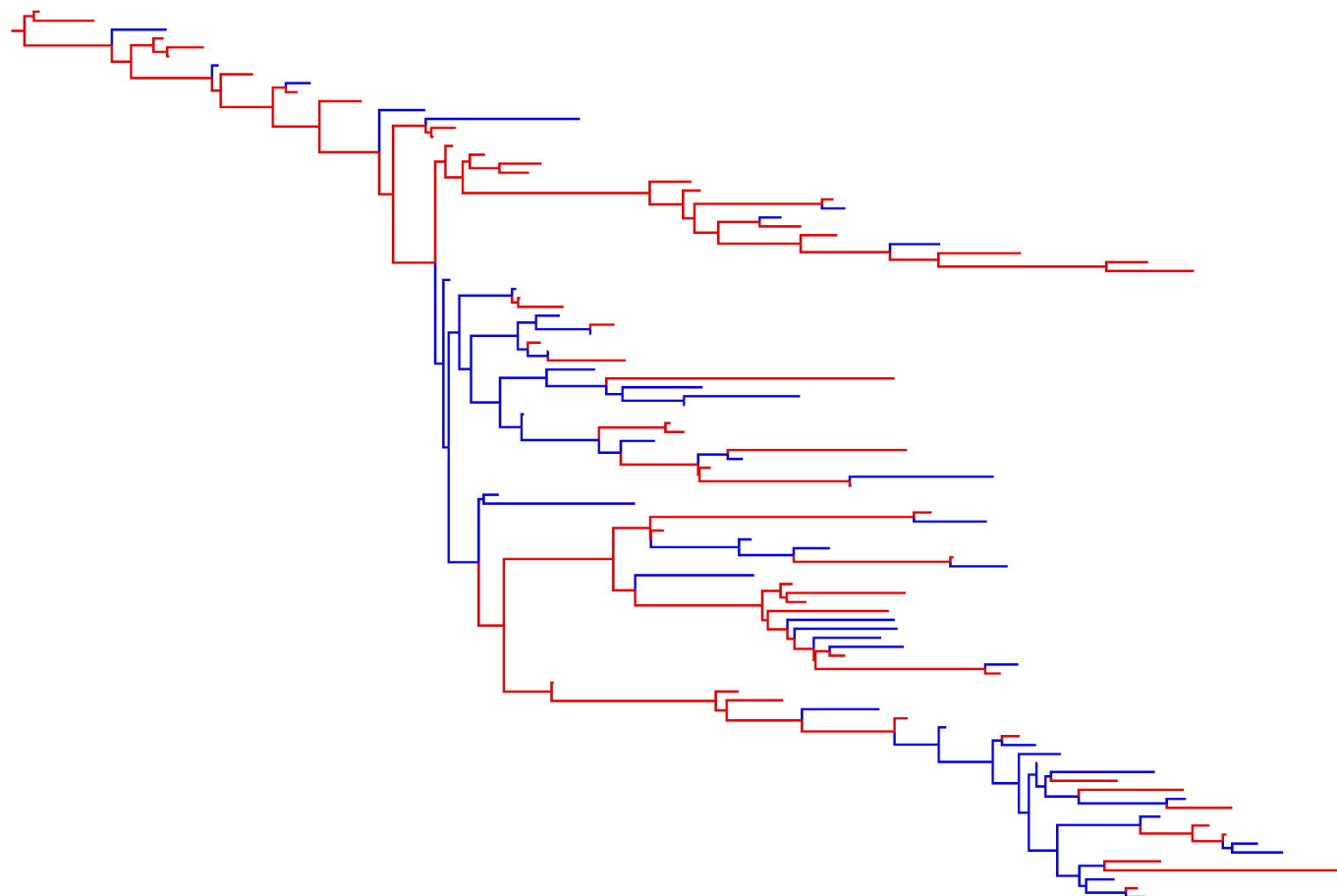
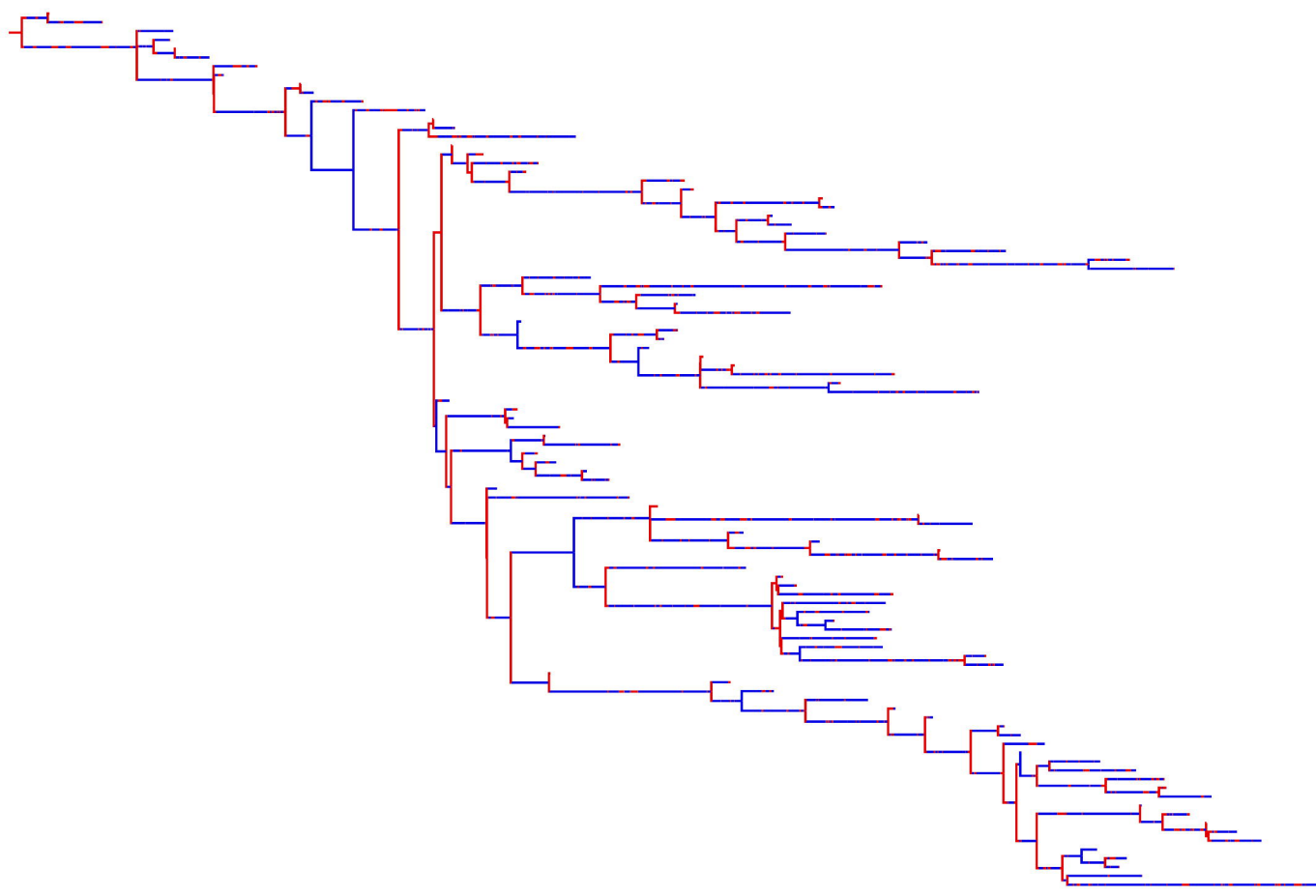
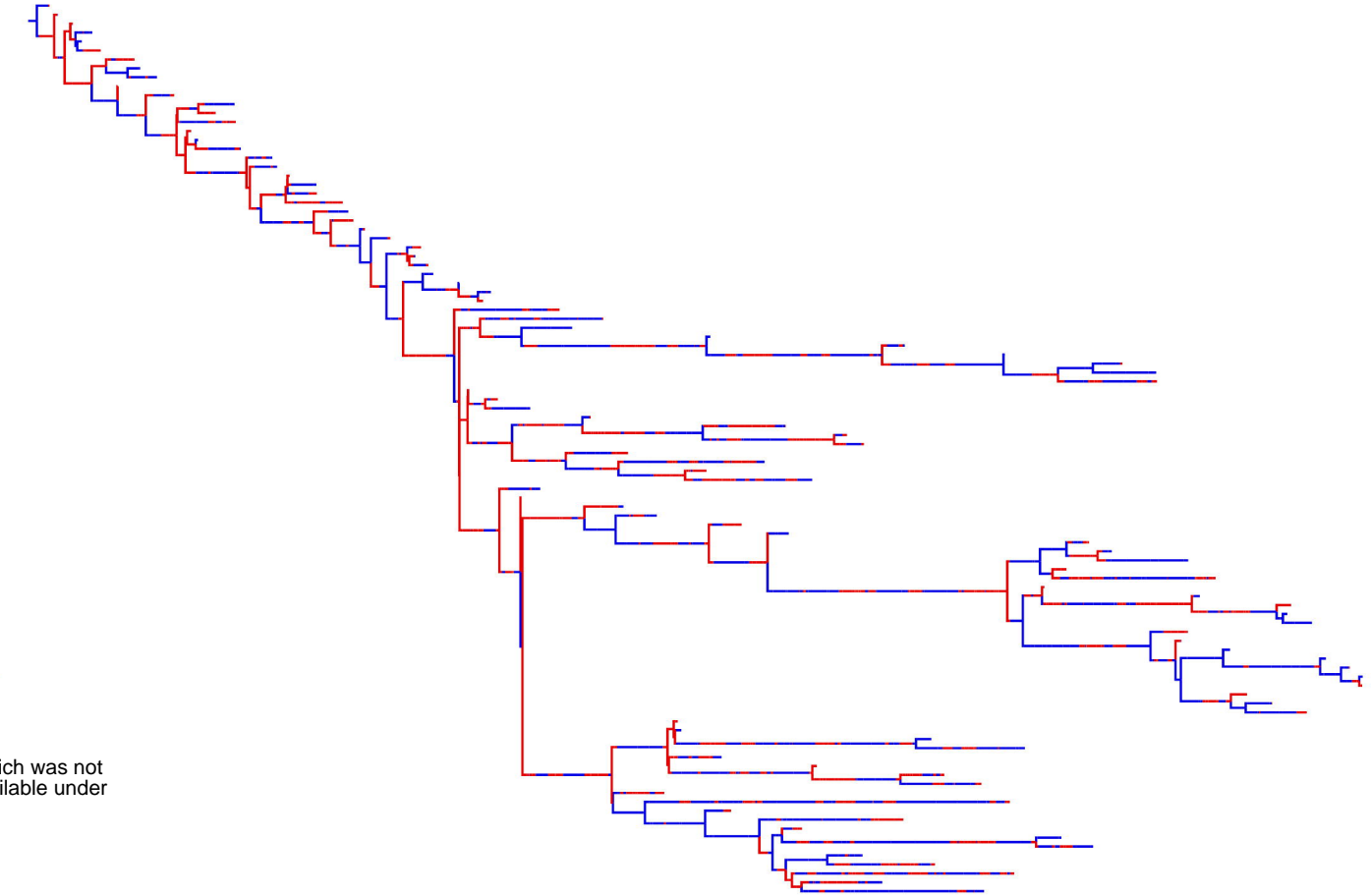
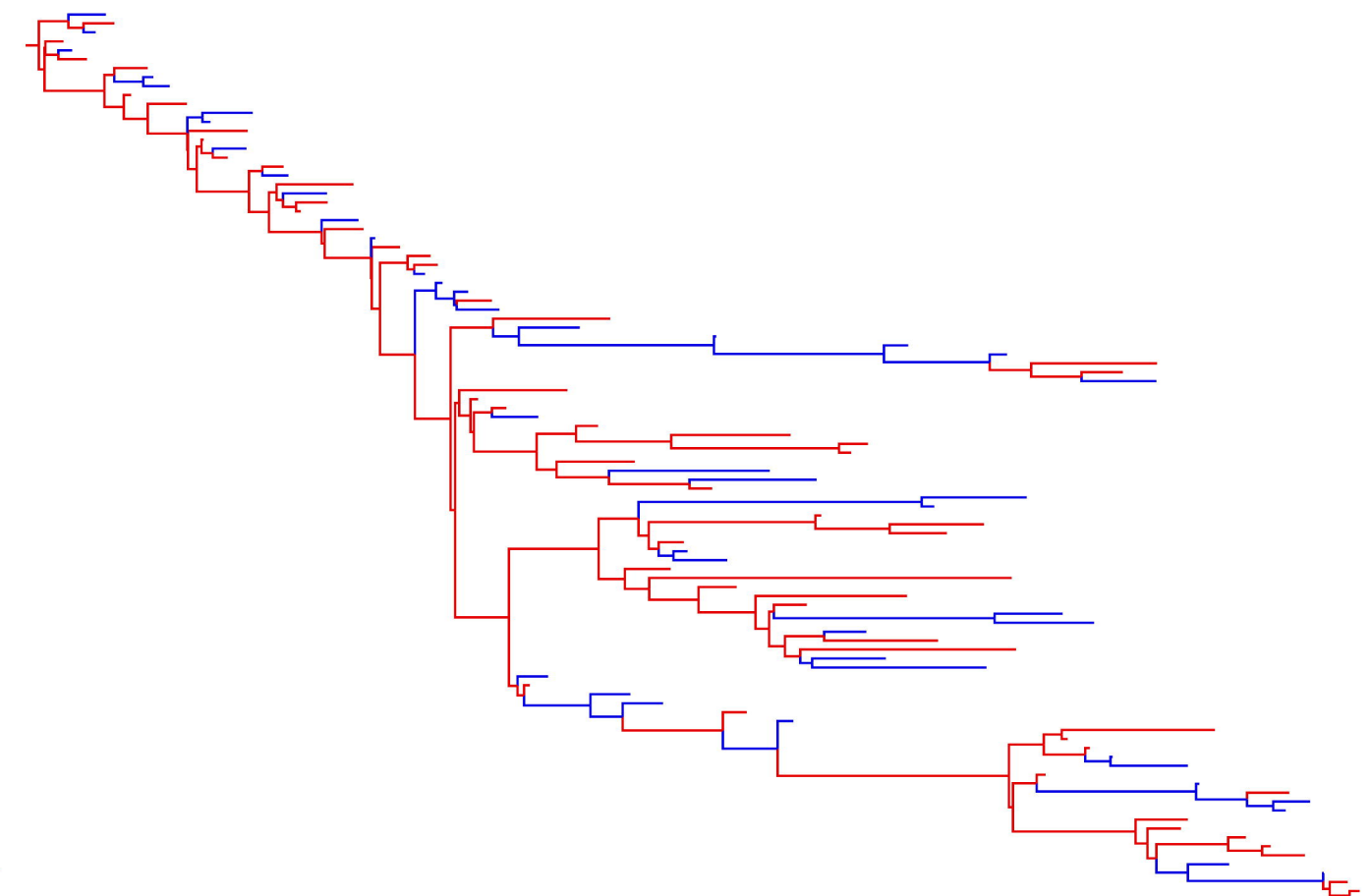
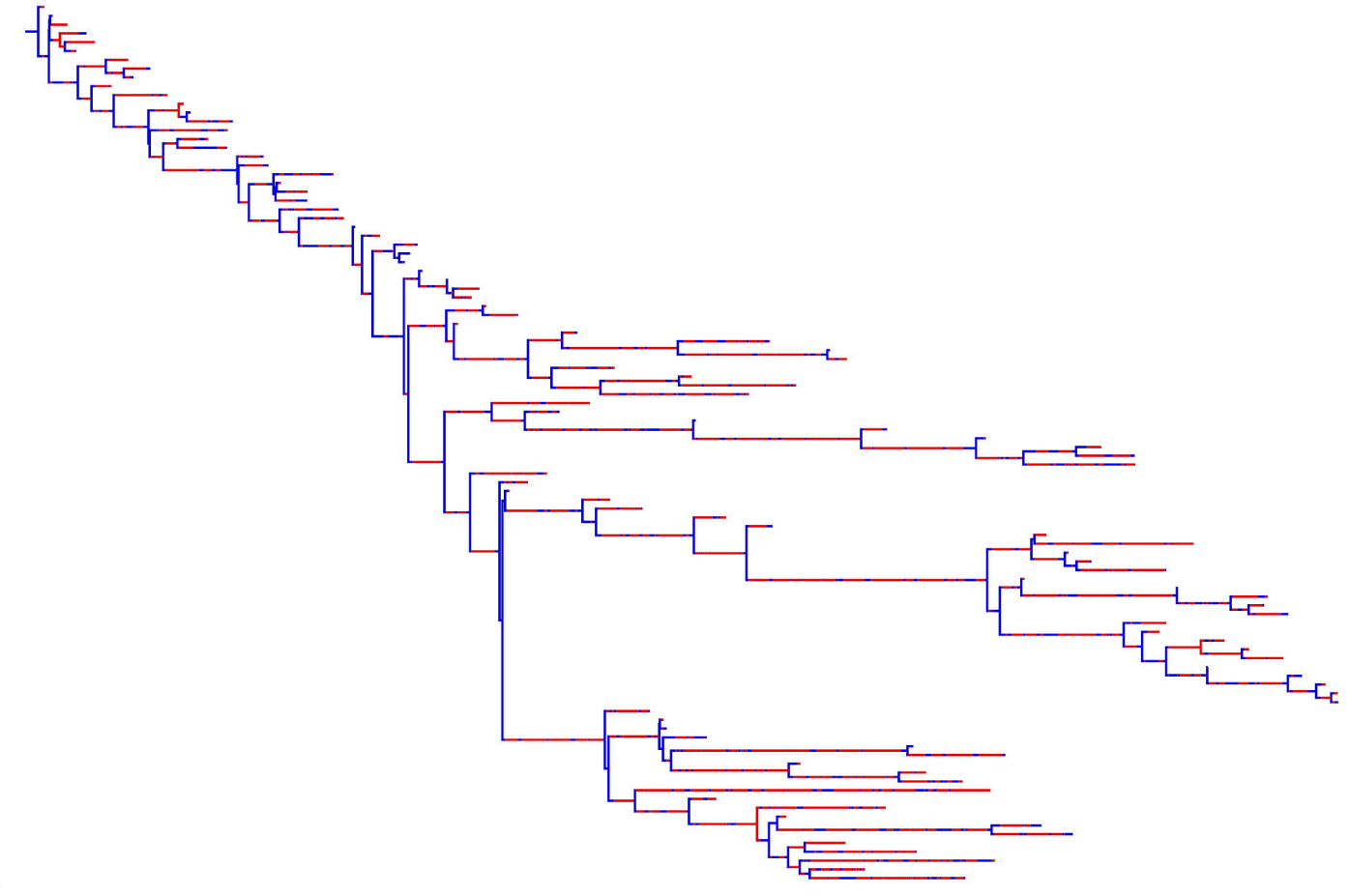


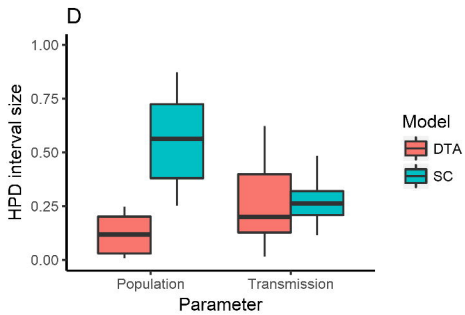
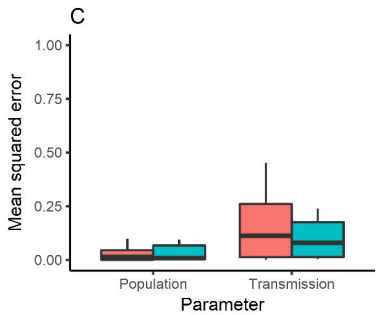
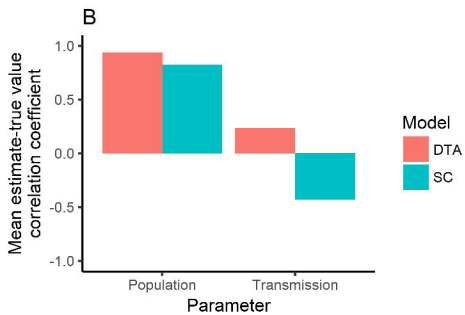
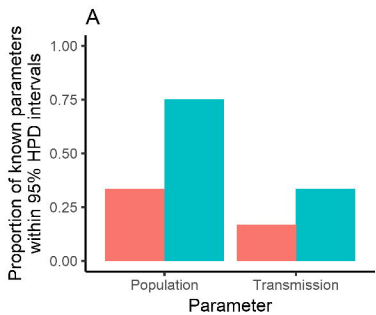
H



A

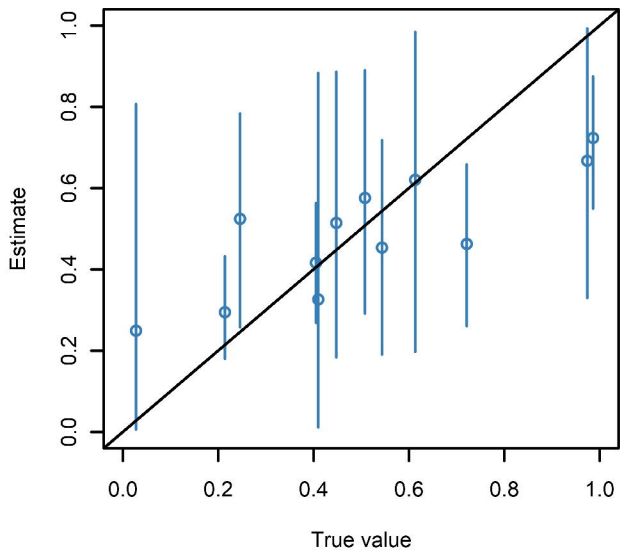
bioRxiv preprint doi: <https://doi.org/10.1101/574087>; this version posted March 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

B**C****D****E****F**

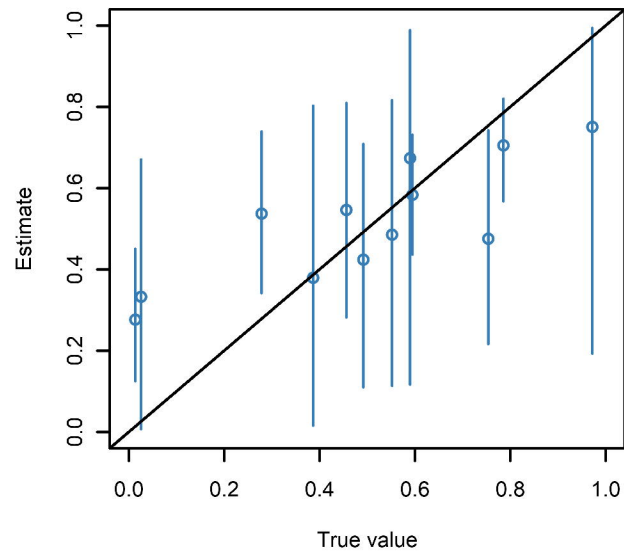


Population

A

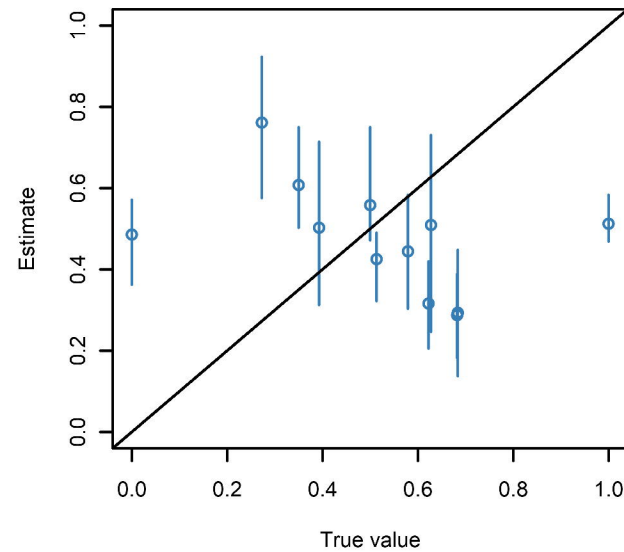


B

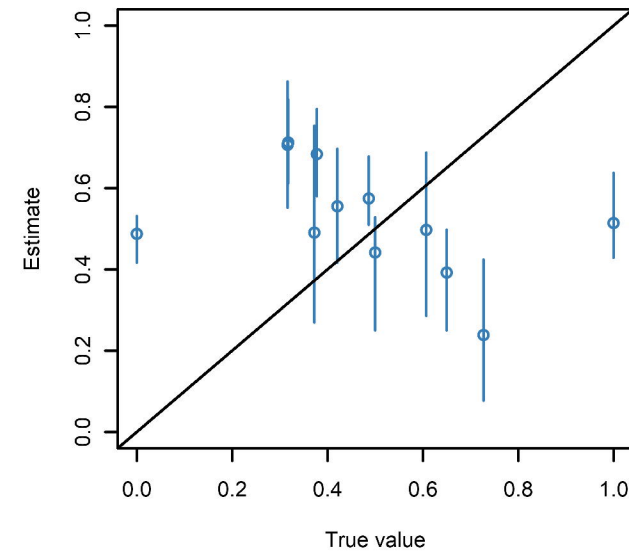


Transmission

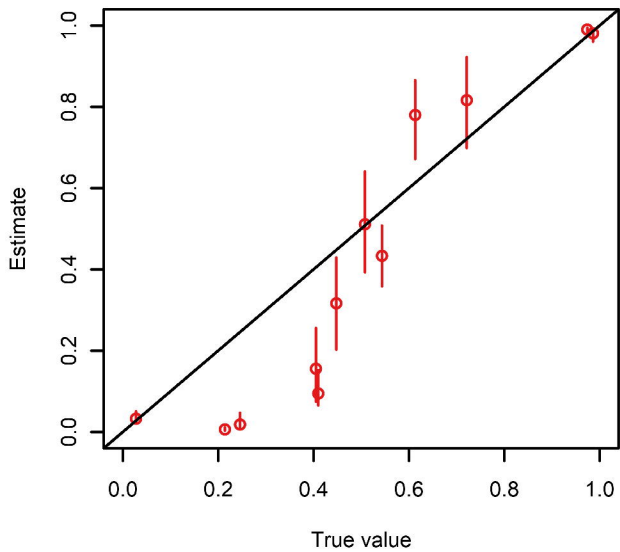
C



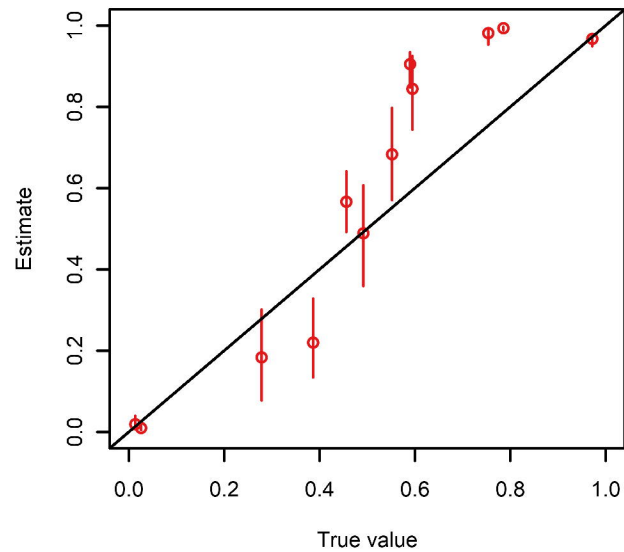
D



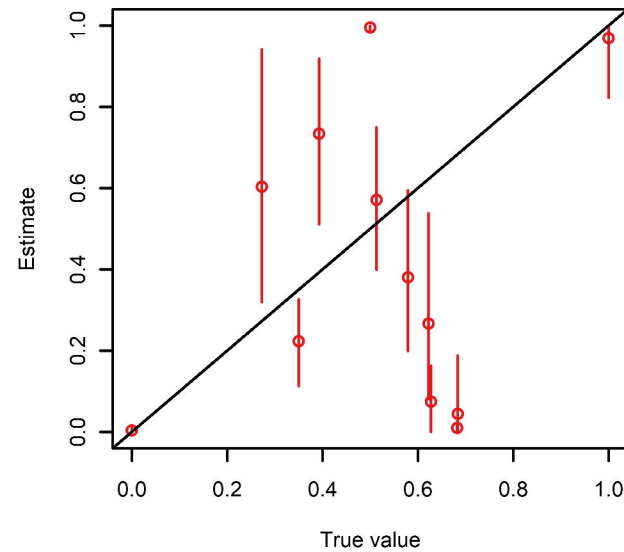
E



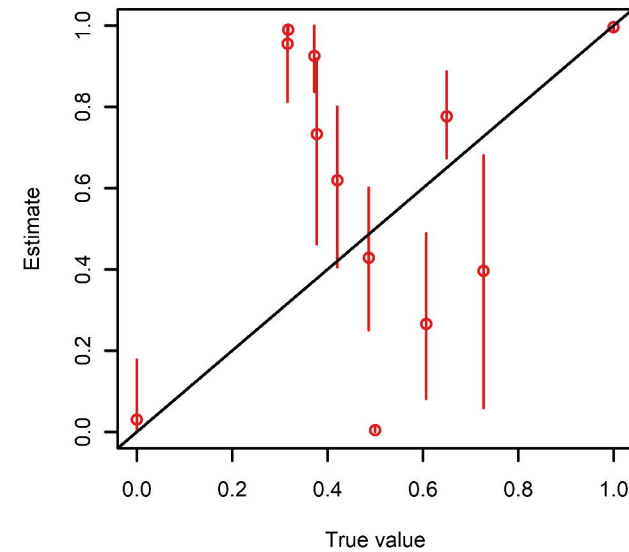
F

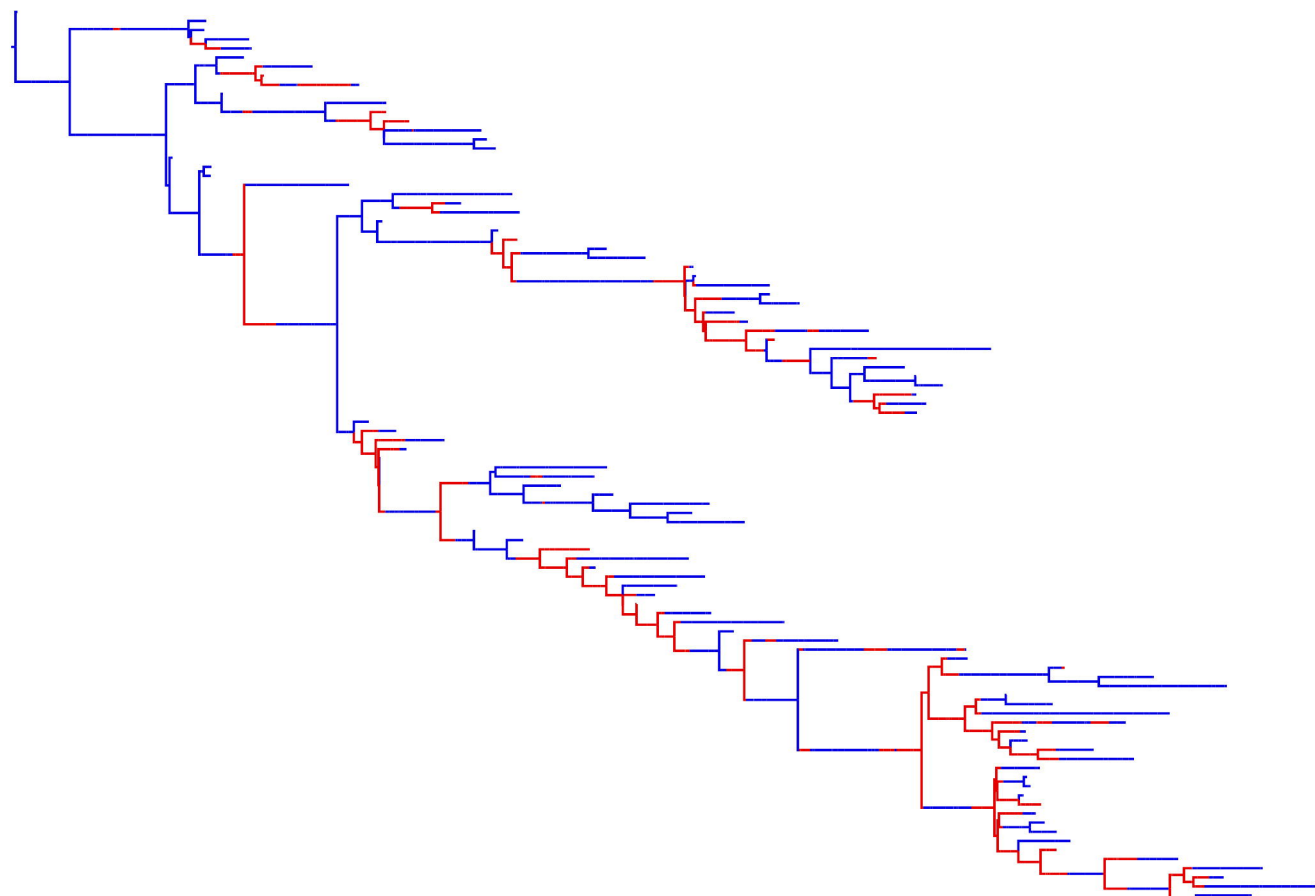


G

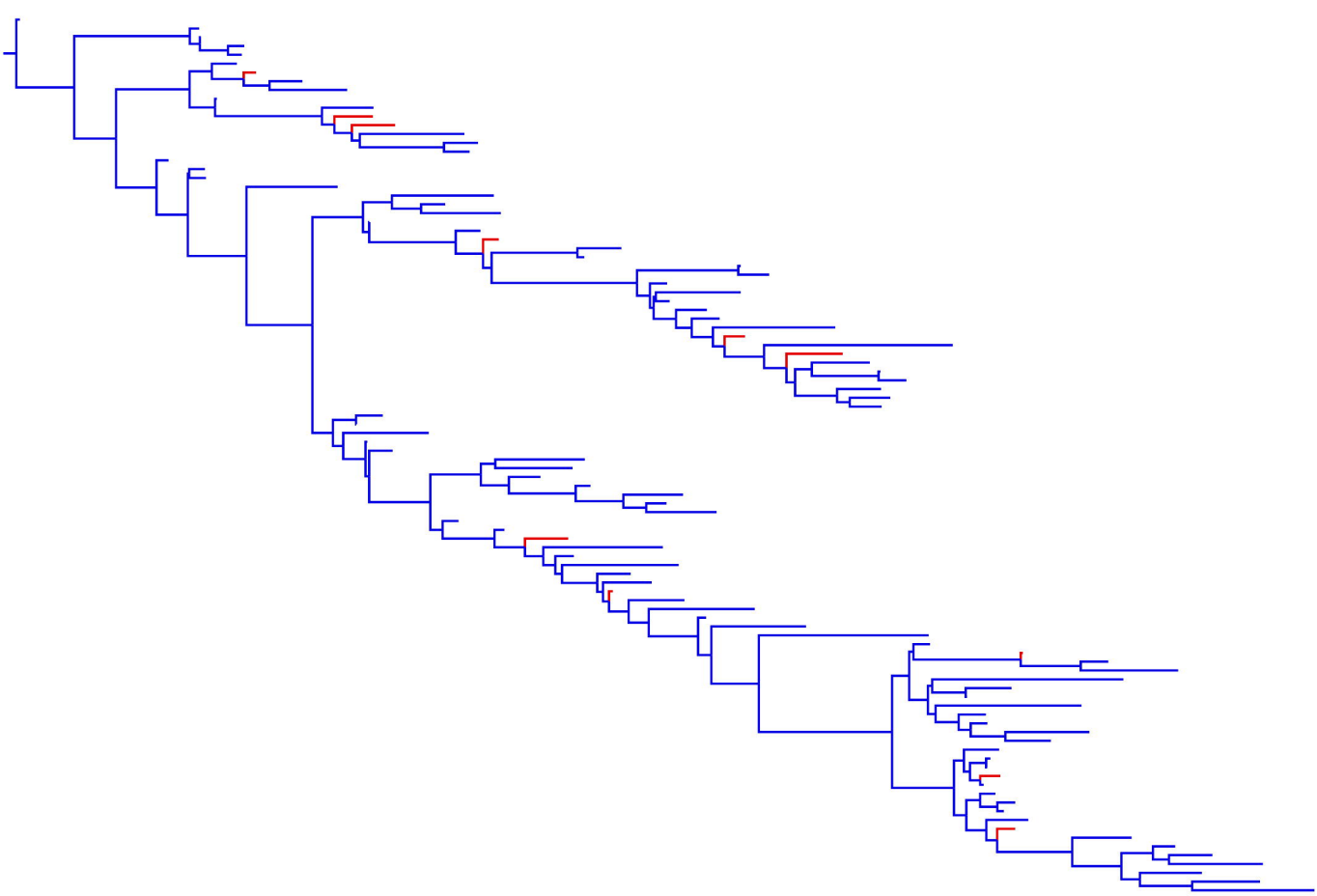
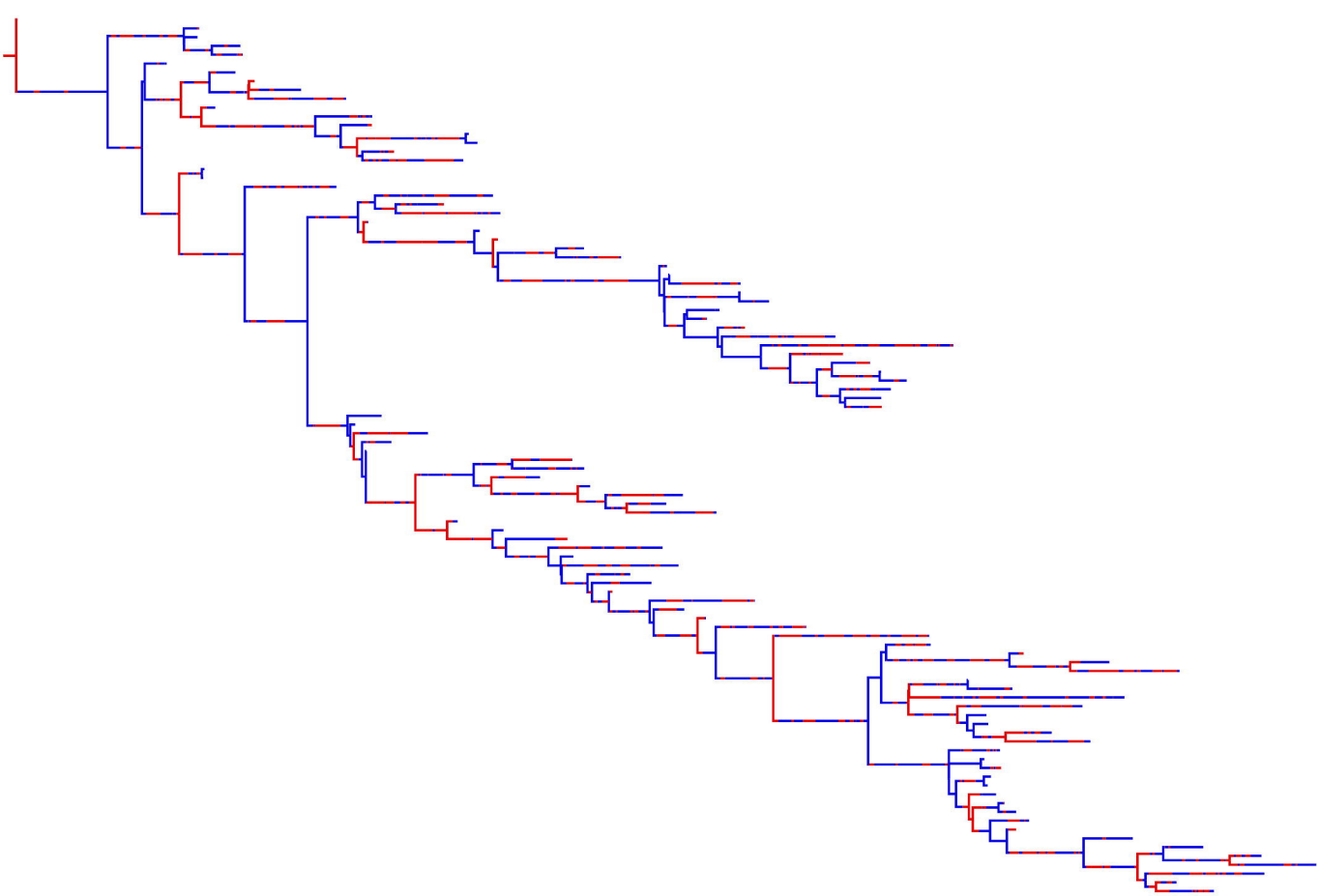


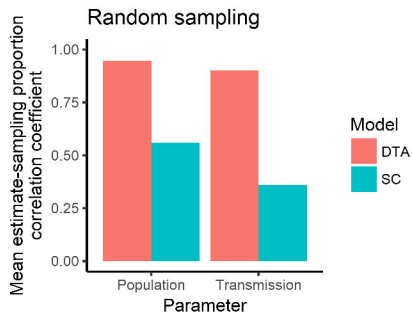
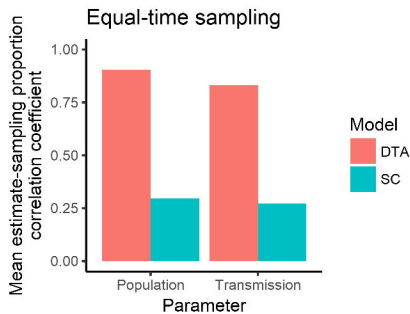
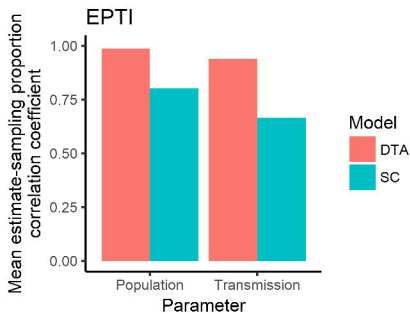
H

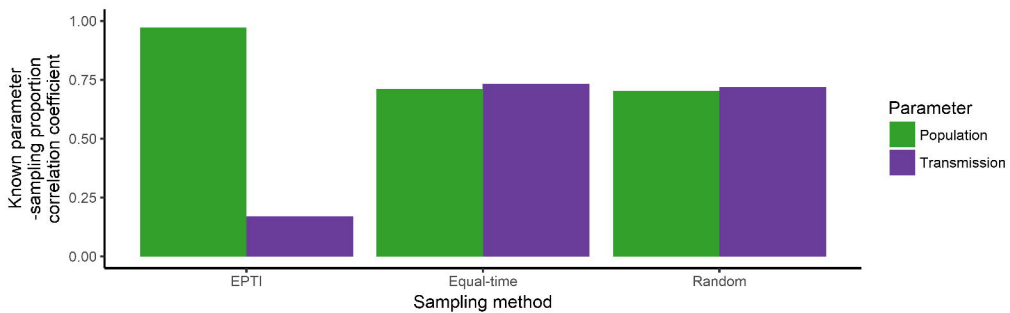


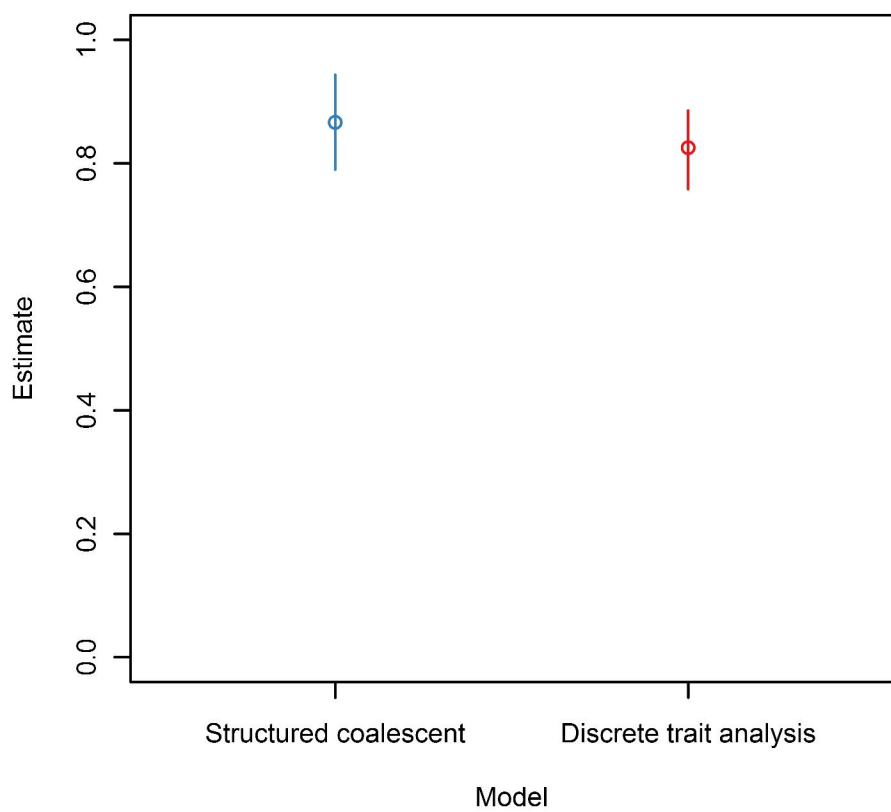
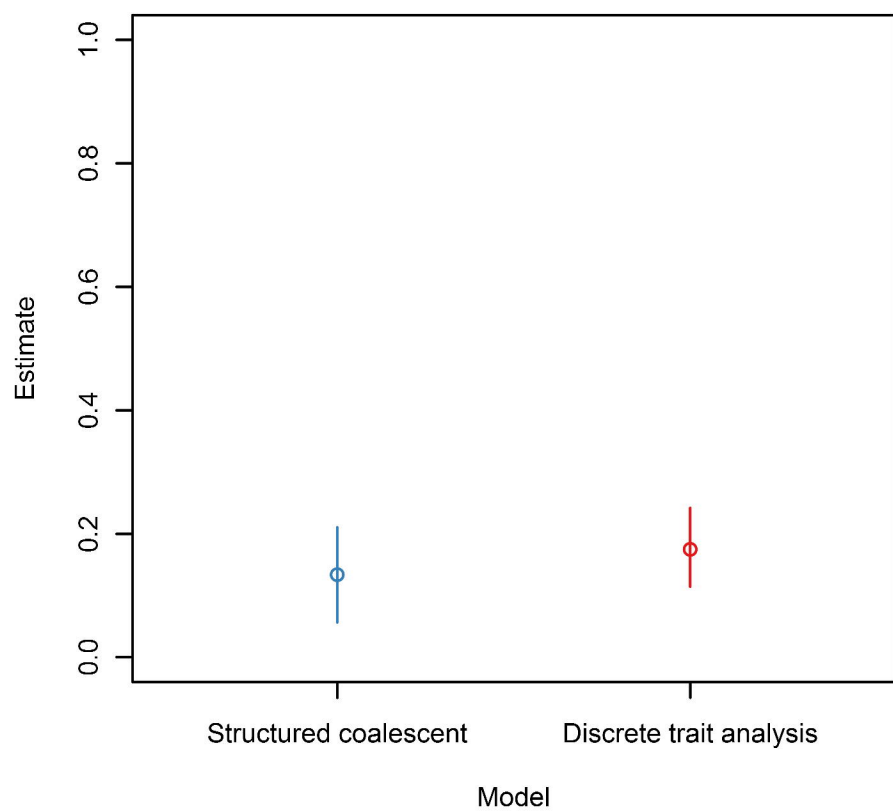
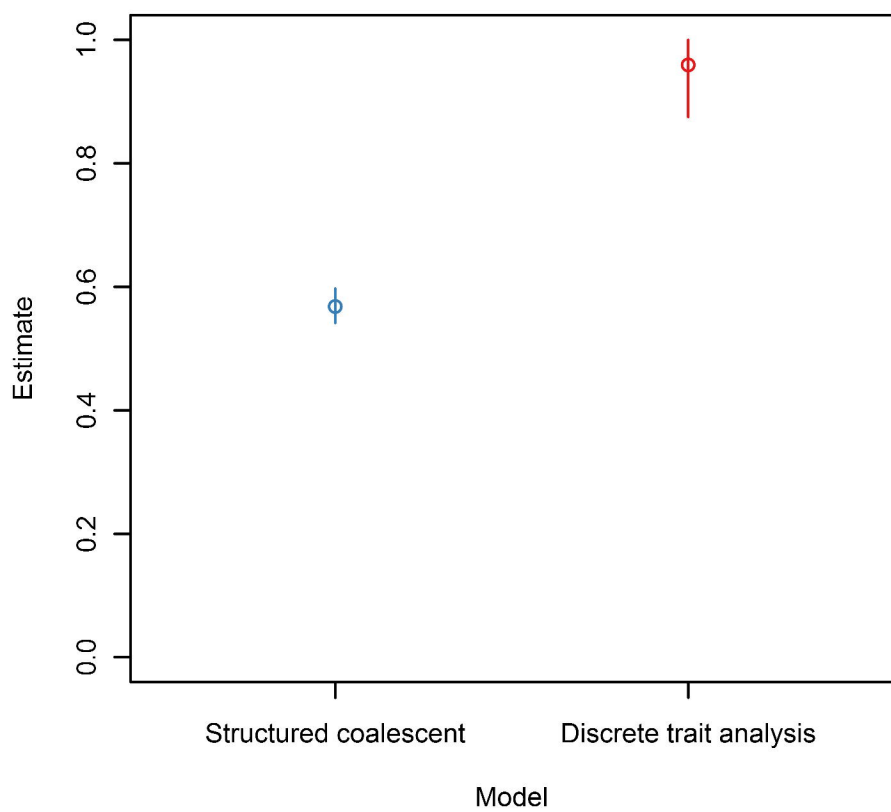
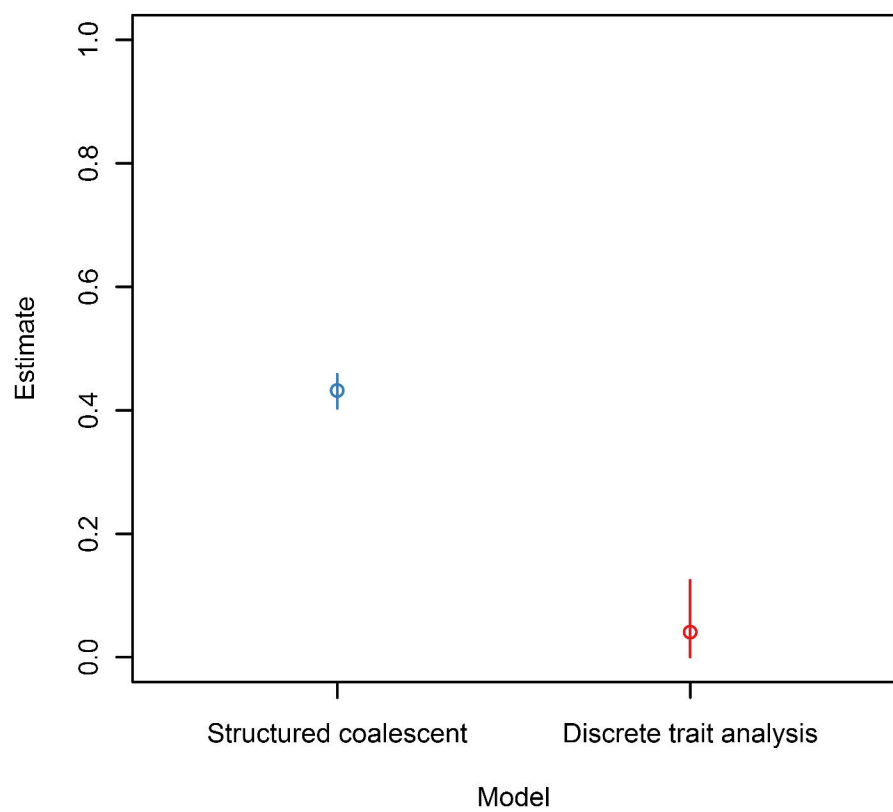
A

bioRxiv preprint doi: <https://doi.org/10.1101/574087>; this version posted March 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

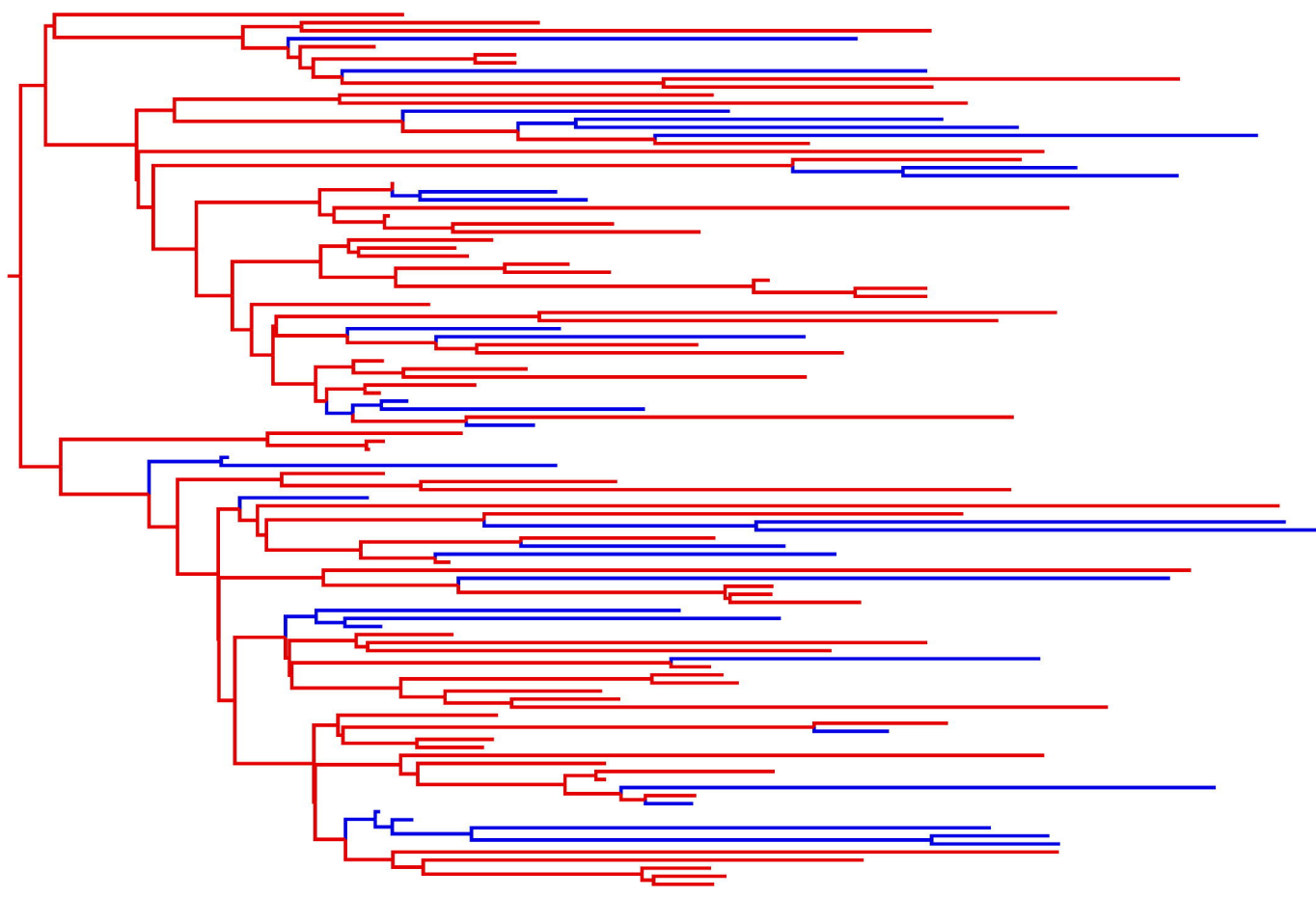
B**C**





A Animal time**B Human time****C Animal-to-human transmissions****D Human-to-animal transmissions**

A



B

