

2019-02-14

1 **The speciation and hybridization history of the genus *Salmonella***

2

3 **Author names**

4 Alexis Criscuolo¹, Sylvie Issenhuth-Jeanjean², Xavier Didelot³, Kaisa Thorell⁴, James Hale⁵,
5 Julian Parkhill⁶, Nicholas R. Thomson⁶, François-Xavier Weill², Daniel Falush^{7,*}, Sylvain
6 Brisse^{8,*}

7

8 **Affiliation**

9 ¹ Bioinformatics and Biostatistics Hub, C3BI, USR 3756 IP CNRS, Paris, France

10 ² Institut Pasteur, Enteric Bacterial Pathogens, World Health Organization Collaborating
11 Centre for Reference and Research on *Salmonella*, Paris, France

12 ³ School of Life Sciences and Department of Statistics, University of Warwick, Coventry,
13 United Kingdom

14 ⁴ Department of Microbiology and Immunology, University of Gothenburg, Gothenburg,
15 Sweden

16 ⁵ Environmental Research Institute, University College Cork, Cork, Ireland

17 ⁶ Wellcome Trust Sanger Institute, Hinxton, United Kingdom

18 ⁷ The Milner Centre for Evolution, University of Bath, Bath, United Kingdom

19 ⁸ Institut Pasteur, Biodiversity and Epidemiology of Bacterial Pathogens, Paris, France

20 ***Corresponding authors**

21 S. Brisse: Biodiversity and Epidemiology of Bacterial Pathogens, Institut Pasteur, 28 rue du
22 Dr Roux, F-75724 Paris, France. E-mail: sylvain.brisse@pasteur.fr; Phone +33 1 45 68 83 34.

23 D. Falush: The Milner Centre for Evolution, Department of Biology and Biochemistry,
24 University of Bath, Bath, United Kingdom. E-mail: danielfalush@gmail.com

25

26 **Keywords:** *Salmonella*; speciation; hybridization; evolution; taxonomy; genomics

2019-02-14

27 **Abstract**

28 Bacteria and archaea make up most of natural diversity but the mechanisms that underlie the
29 origin and maintenance of prokaryotic species are poorly understood. We investigated the
30 speciation history of the genus *Salmonella*, an ecologically diverse bacterial lineage, within
31 which *S. enterica* subsp. *enterica* is responsible for important human food-borne infections.
32 We performed a survey of diversity across a large reference collection using multilocus
33 sequence typing, followed by genome sequencing of distinct lineages. We identified eleven
34 distinct phylogroups, three of which were previously undescribed. Strains assigned to
35 *S. enterica* subsp. *salamae* are polyphyletic, with two distinct lineages that we designate
36 Salamae A and Salamae B. Strains of subspecies *houtenae* are subdivided into two groups,
37 Houtenae A and B and are both related to Selander's group VII. A phylogroup we designate
38 VIII was previously unknown. A simple binary fission model of speciation cannot explain
39 observed patterns of sequence diversity. In the recent past, there have been large scale
40 hybridization events involving an unsampled ancestral lineage and three distantly related
41 lineages of the genus that have given rise to Houtenae A, Houtenae B and VII. We found no
42 evidence for ongoing hybridization in the other eight lineages but detected more subtle signals
43 of ancient recombination events. We are unable to fully resolve the speciation history of the
44 genus, which might have involved additional speciation-by-hybridization or multi-way
45 speciation events. Our results imply that traditional models of speciation by binary fission and
46 divergence may not apply in *Salmonella*.

47

2019-02-14

48 **Data summary**

49 Illumina sequence data were submitted to the European Nucleotide Archive under project
50 number PRJEB2099 and are available from INSDC (NCBI/ENA/DDBJ) under accession
51 numbers ERS011101 to ERS011146. The MLST sequence and profile data generated in this
52 study have been publicly available on the *Salmonella* MLST web site between 2010 and the
53 migration of the *Salmonella* MLST website to EnteroBase
54 (<https://enterobase.warwick.ac.uk/>), and subsequently from there.

55

2019-02-14

56 **Introduction**

57 Bacteria and archaea make up most of natural diversity, both in terms of species richness and
58 biological functions [1,2]. However, the mechanisms that underlie the origin and maintenance
59 of prokaryotic species are poorly understood. It is often assumed that there is a single
60 phylogenetic tree representing the relationships amongst prokaryotic taxa, with the branch
61 lengths reflecting divergence times between them. However, bacteria and archaea acquire
62 foreign DNA by homologous and non-homologous recombination and can recombine
63 frequently, including in the *Salmonella* genus [3–9]. High recombination rates can maintain
64 genetic cohesion within a species, preventing divergence and speciation from occurring until
65 barriers to gene flow develop. Recombination has been shown in laboratory experiments to be
66 suppressed by nucleotide mismatches between donor and recipient [10,11]. This property
67 provides a potential mechanism for speciation. It has been shown by simulation that large
68 effective population sizes and neutral genetic drift can precipitate speciation by increasing the
69 average pairwise divergence between strains, leading to either binary or multi-way speciation
70 events [5,12].

71 Conversely, distinct new lineages or species can potentially arise almost instantaneously by
72 hybridization of existing distantly related ones. Such large-scale hybridization events can
73 occur at once by recombination of large genomic regions (e.g., [13]), or through multiple
74 exchanges of small chromosomal segments associated with ecological convergence [14].
75 Therefore, to describe relationships between prokaryotes and understand patterns of species
76 richness and phenotypic diversity, it is important to characterise the process of speciation and
77 gene flow between species, including large-scale hybridization events [15].

78 Salmonellae are a prominent speciation model, where experimental and genomic studies of
79 recombination and hybridization have been pioneered [4-10,14]. The genus *Salmonella* is
80 divided into a number of phylogroups, namely *bongori*, *enterica*, *salamae*, *arizonae*,
81 *diarizonae*, *houtenae*, and *indica* [16–18]. *Salmonella bongori* has been classified a distinct
82 species [18], while the other phylogroups are considered to be subspecies of a single species,
83 *S. enterica*. These taxa are further subdivided into serovars based on antigenic variation of
84 flagellins and O-antigen.

85 Members of the genus *Salmonella* are major pathogens of humans and other warm-blooded
86 animals. Human infections mostly involve *S. enterica* subspecies *enterica*, which can cause
87 gastroenteritis, enteric fever and other infections [19,20]. Other *S. enterica* subspecies, as well

2019-02-14

88 as the species *S. bongori*, are more typically isolated from cold blooded animals or the
89 environment, and are rarely reported from human infections [21].

90 Here we are concerned with evolutionary relationships rather than taxonomy and we
91 designate phylogroups by names that derived from these subspecies' labels, e.g. Bongori,
92 Arizonae, Diarizonae, etc., with Enterica representing subspecies *enterica*. We use italicised
93 names such as *houtenae* to refer to previous subspecies designations, which sometimes differ
94 from our phylogroup assignments. A seventh *S. enterica* subgroup (group VII) was
95 distinguished based on multilocus enzyme electrophoresis and gene sequencing [22–24]. Note
96 that phylogenetic re-evaluation [25] of the proposed species *Salmonella subterranea* [26]
97 shows that it does not belong to the *Salmonella* genus.

98 Phylogenetic analyses of the evolutionary relationships amongst the different *Salmonella*
99 lineages have led to contradictory conclusions with several proposed phylogenetic trees
100 [9,23,24,27–37]. This lack of consensus might reflect technical issues with phylogenetic
101 reconstruction but a more biologically interesting possibility is that the history of *Salmonella*
102 is not well-characterized by a simple model in which speciation proceeds stepwise by
103 irreversible binary fissions.

104 To test this hypothesis, we sampled the genetic diversity within the little studied groups from
105 cold-blooded hosts and used whole genome sequences from representative isolates of
106 phylogroups to characterize the genetic relationships between them and to infer historical
107 populations splits and gene flow. We show that while a binary fission model of speciation
108 works for some of the *Salmonella* lineages, there are several important historical events that
109 cannot be characterized in this way.

2019-02-14

110 **Methods**

111 *Taxonomic sampling and MLST analyses*

112 A total of 367 strains (**Table S1**) from outside the subspecies enterica were selected from the
113 collection of the World Health Organization Collaborative Centre for Reference and Research
114 on *Salmonella*, Institut Pasteur, Paris, France. This center contains the reference strains of all
115 *Salmonella* serovars and their variants. The 367 strains represented approximately one third of
116 currently described serovars outside enterica and were selected to maximize the diversity of
117 antigenic formulae. MLST was performed on these strains using updated primers adapted
118 from those of Kidgell et al. [38] to amplify DNA from *S. bongori* and all subspecies of
119 *S. enterica*. The novel primers are described in **Table S3**; note that they have been publicly
120 available on the MLST web site between 2008 and the migration of the *Salmonella* MLST
121 website to EnteroBase, and subsequently from there.

122 A phylogenetic tree was inferred from the median distance matrix of the seven genes with the
123 algorithm BioNJ* [39]. A supermatrix of characters was built by concatenating the seven MSAs
124 with the program Concatenate (www.supertriplets.univ-montp2.fr/PhyloTools.php), and the
125 nucleotide diversity of groups was defined using the index π [40] with the program DnaSP
126 [41]. Minimum spanning trees were built using the software tool BioNumerics (Applied-
127 Maths, Belgium).

128

129 *Strain selection and genome sequencing*

130 A set of 46 strains were selected for whole genome sequencing (**Table S2**). Genome
131 sequencing was achieved by Illumina 2 x 50 nt paired-end sequencing for all strains. The
132 characteristics of the obtained *de novo* assemblies are summarised in **Table S2**. This set was
133 completed with genome sequences gathered from the GenBank repository (*i.e.*, 16 *S. enterica*
134 *subsp. enterica*, 1 *S. enterica subsp. arizonae*, and 1 *S. bongori* strains), as well as 9
135 *S. bongori* genome sequences from Fookes et al. [34]. This led to a total of 73 genomes
136 (**Table S2**): *S. enterica subsp. enterica*, 16; *subsp. salamae*, 13; *subsp. arizonae*, 9; *subsp.*
137 *diarizonae*, 10; *subsp. houtenae*, 6; *subsp. indica*, 4; *S. bongori*, 10; and VII, 2.

138

139 *Core gene construction*

2019-02-14

140 Each of the 4,423 protein sequences from *S. enterica* strain LT2 [42] was used as query to
141 perform BLAST similarity searches [43] against the genome sequence of each of the other 72
142 strains. Clusters of homologous sequences were built by considering only the first tblastn hit
143 (E-value < 10^{-5}), and every cluster that did not contain 73 sequences (*i.e.*, one per strain) was
144 discarded. Next, orthology was assessed within each cluster by performing reciprocal tblastn,
145 leading to 2,328 clusters of putative orthologous coding sequences from the core gene set of
146 the 73 strains. For each of these clusters, sequences were translated, and a multiple sequence
147 alignment (MSA) was performed with ProbCons [41] and next back-translated to obtain a
148 codon-level MSA. The 2,328 MSAs were concatenated into a supermatrix of 2,137,446
149 nucleotide characters that was used to infer a balanced minimum-evolution phylogenetic tree
150 using FastME [44] with pairwise p-distances. Branch support was assessed for each internal
151 branch with a MSA-based bootstrap procedure based on 1,000 replicates. This procedure
152 samples the MSA with replacement according to the same procedure as the standard bootstrap
153 with nucleotide characters.

154

155 *Recombination analysis*

156 We applied four separate and complementary methods to analyse the ancestral recombination
157 events that occurred during the evolution of the genus *Salmonella*. Firstly, we applied
158 chromosome painting on the 73 genomes, using CHROMOPAINTER [45] to reconstruct each
159 genome as a mosaic of all the others. The results were summarized as a heatmap of coancestry
160 values, where each coancestry value corresponds to the number of fragments copied from one
161 genome to another (**Figure 2**). Secondly, we performed pairwise comparisons of genomes
162 using a gene-by-gene approach. For each pair of genomes, we computed the genetic distances
163 of all shared genes, and the distribution of these distances was plotted as a cumulative curve
164 (**Figure 3**). Thirdly, the CHROMOPAINTER analysis was repeated using only nine unrelated
165 genomes: one for each of the 12 phylogroups but excluding VII and Houtenae B due to recent
166 shared ancestry with Houtenae A, and considering Enterica A and B as a single group. Each
167 genome was therefore reconstructed as a mosaic of the other eight unrelated genomes. This
168 allowed us to explore deeper relationships between phylogroups, since when all genomes are
169 included each genome from a phylogroup copies mostly from other genomes of the same
170 phylogroup (**Figure 2**). The resulting coancestry matrix was plotted as a heatmap (**Figure 4**).
171 Fourthly, we applied the Treemix algorithm with parameter K=3 [46] to one genome from

2019-02-14

172 each of the 12 phylogroups in order to reconstruct their relationships as a vertical
173 phylogenetic tree augmented with horizontal transfer events (**Figure 5**).

174

175 *Pan-genome analyses*

176 Analysis of accessory genome was performed using the Roary pan-genome pipeline version
177 3.6.2 [47]. Since the draft genomes were very unequally fragmented and synteny information
178 therefore was of variable reliability we used the “don't-split-paralogs” option. The analysis
179 was performed with a protein identity cut-off of 85% and the core genome was defined as
180 genes present in > 99% of the genomes studied. The Pearson correlation between accessory
181 gene content of the genomes were visualised using the R software CORRPLOT package [48].

182

2019-02-14

183 **Results and Discussion**

184 In order to survey the diversity of *Salmonella* outside *S. enterica* subsp. *enterica*, a total of
185 367 strains, comprising about a third of the known non-Enterica serovars, were selected from
186 the World Health Organization Collaborating Centre for Reference and Research on
187 *Salmonella* (Institut Pasteur, Paris, France) reference collection and subjected to multilocus
188 sequence typing (MLST) (**Tables S1, S2**). A phylogenetic tree was built (**Figure S1**),
189 revealing a novel group (labelled VIII) and suggesting a polyphyletic origin of *S. enterica*
190 subsp. *salamae* (Salamae A and B) and of *S. enterica* subsp. *houtenae* (Houtenae A and B).
191 Within-phylogroup nucleotide diversity (**Figure S1 inset**) was the highest in Arizonae
192 ($\pi = 1.6\%$), lowest in Houtenae groups, Bongori, Salamae B and Diarizonae (π ranging from
193 0.35% to 0.42%), whereas it was intermediate in Salamae A, Indica and Enterica. Minimum
194 spanning tree analysis of MLST profiles illustrates the genotypic diversity within each group
195 (**Figure S2**).

196 Based on MLST diversity, 46 genomes were chosen for genome sequencing and compared to
197 27 previously published genome sequences of Enterica, Arizonae, and Bongori (**Table S2**). A
198 phylogenetic tree based on the genome sequences is shown in **Figure 1**. This tree implies that
199 *S. enterica* subsp. *salamae* is not a monophyletic group but instead forms two lineages with
200 distinct evolutionary histories that we designate Salamae A and Salamae B. Whereas Salamae
201 A contained 138 (88%) of the *salamae* strains, Salamae B comprised 18 isolates collected
202 from a human (one isolate), a bat (one isolate) or reptiles (16 isolates, including 6 from
203 chameleons). In contrast, 49 (41.5%) of Salamae A isolates were from humans, and only 34
204 (28.8%) were from cold-blooded animals, suggesting important ecological and pathogenic
205 differences between the two Salamae groups. *S. enterica* subsp. *houtenae* was also subdivided
206 into two distinct phylogroups, which we have designated Houtenae A and Houtenae B, and
207 which clustered together with group VII on the tree. The genome-wide phylogenetic analysis
208 also uncovers a hitherto unknown phylogroup, labelled VIII, made of strains previously
209 identified as either *salamae*, *diarizonae* or of the former Hisingen serotype of *S. enterica*
210 subsp. *enterica* [25]. The description of Salamae B, Houtenae B and VIII represent the first
211 novel *Salmonella* phylogroups described since the identification of group VII by Selander and
212 colleagues more than 25 years ago [24,27]. Our analysis therefore defines 11 phylogroups
213 within *Salmonella*. The phylogenetic tree also shows further subdivisions at shallower levels,
214 including the division of *S. enterica* subsp. *enterica* into Enterica A and Enterica B as
215 previously described [5,9]. Note that the genomes of the present study have been publicly

2019-02-14

216 available from International Nucleotide Sequence Database Collaboration (INSDC) since
217 2011, and were used in a genome-based phylogenetic analysis of *Salmonella* by Alikhan *et al.*
218 [49]; the three novel *Salmonella* groups were labelled as novel subspecies A (Houtenae B), B
219 (VIII) and C (Salamae B) in [49].

220

221 **Recent recombination between phylogroups.**

222 We used chromosome painting of the above set of 73 strains in order to investigate shared
223 ancestry and recombination events between different phylogroups. Specifically, the
224 CHROMOPAINTER algorithm uses a Hidden Markov Model to reconstruct each isolate as a
225 mosaic of stretches of DNA of the other isolates in the sample [45]. The results can be
226 summarized as a heatmap indicating how many stretches from each other sample are used in
227 the reconstruction. The organism used in the reconstruction is assumed to be the most closely
228 related for each stretch of DNA. **Figure 2** shows a heatmap illustrating the proportion of
229 DNA used to paint each isolate across the genome, with dark blue corresponding to 0% and
230 dark red corresponding to 10%. We call this proportion the coancestry value. Each
231 phylogroup shows higher coancestry within the same phylogroup than with others. The
232 highest coancestry between strains in different phylogroups is between Houtenae A, Houtenae
233 B and VII. However, Houtenae B shows higher Enterica ancestry (particularly with Enterica
234 B) than do Houtenae A or VII. The two deepest branching Salamae A strains show high levels
235 of coancestry with several other groups including Salamae B, Diarizonae, Indica and VIII.
236 One strain of Enterica A (SL483) is exceptional in showing higher coancestry levels with
237 Enterica B.

238 In order to test whether high coancestry between groups might be explained by recent
239 recombination between them, we looked for evidence of sharing of very similar stretches of
240 DNA between pairs of lineages [14] by plotting, for each pairwise comparison, the proportion
241 of genes with divergence below a threshold increasing from 0% to 25% (**Figure 3**).
242 Consistent with recent recombination between them, Enterica B and Houtenae B showed
243 many more genes with very similar sequences than expected based on their position in the
244 phylogenetic tree, with 20% of the genes of an Enterica B strain having less than 1%
245 divergence to Houtenae B, compared to only 5% between Enterica B and Houtenae A (**Figure**
246 **3A**). These divergence curves are also consistent with recent recombination between
247 Houtenae A, Houtenae B and VII. For example, approximately 5% of the VII genome and 6%
248 of Houtenae A has less than 0.1% divergence with Houtenae B (**Figure 3B**), suggesting that

2019-02-14

249 there has been very recent recombination between these three phylogroups. There is no
250 analogous signal of recent recombination between any of the strains of Salamae A or Salamae
251 B with each other or with other phylogroups based on cumulative divergence curves (*e.g.*,
252 **Figure 3C**). The smudged pattern of coancestry of the deeper branching Salamae A and
253 Salamae B strains in **Figure 2** can potentially be explained by them retaining ancestral
254 variants that have been lost by the rest of the phylogroup and therefore does not necessarily
255 indicate recent recombination between lineages. **Figure 3D** illustrates the absence of any
256 signal of recent recombination with Arizonae.

257

258 **Evidence for hybridization in the origin of the phylogroups**

259 We next examined the origins of the phylogroups themselves. Recombination events which
260 predate the generation of the diversity observed *within* each phylogroup are unlikely to be
261 picked up in the chromosome painting analysis in **Figure 2**: members of a phylogroup that
262 have inherited the same ancestrally imported stretch will be painted by each other for those
263 stretches. Therefore, we selected a single strain from each phylogroup and performed a
264 distinct chromosome painting analysis. We excluded VII and Houtenae B due to the recent
265 shared ancestry with Houtenae A, and also included only a single representative for both
266 Enterica A and Enterica B. The chromosome painting results (**Figure 4**) show high coancestry
267 between Bongori and Arizonae and between Indica and Enterica. These relationships can be
268 interpreted using a vertical phylogenetic model, as they agree with a large number of different
269 analyses including ours (**Figure 1**) that Arizonae is the deepest branching lineage after
270 Bongori and that Indica is a sister group of Enterica [9,24,33,34]. Conversely, the
271 chromosome painting analysis reveals a large number of intransitive relationships (*i.e.*, in
272 which A has elevated coancestry with B and B has high coancestry with C but C does not
273 have high coancestry with A). First, Diarizonae and Arizonae have high coancestry, as do
274 Diarizonae and Salamae B but Salamae B and Arizonae do not (**Figure 4**). Second, Houtenae
275 A and Salamae A have high coancestry with each other and the phylogenetic tree suggests
276 they are sister taxa. However, they have different relationships to other phylogroups.
277 Houtenae A, but not Salamae A, shows high coancestry with Arizonae, while Salamae A
278 shows higher shared ancestry with Indica and Enterica. Intransitive patterns of coancestry are
279 also evident for VIII, Salamae B and Diarizonae and for VIII, Salamae B and Bongori. An
280 intransitive pattern is not predicted by any phylogenetic model and is indicative of mixture in

2019-02-14

281 the history. These observations suggest a complex pattern of homologous recombination
282 events that predate diversification within phylogroups.

283

284 **A scenario involving three recent hybridization events**

285 To complement the results above, we used Treemix to infer a history that allows for
286 recombination events in the origins of the phylogroups. Treemix attempts to model the
287 covariance matrix reflecting SNP sharing between strains by assuming a phylogenetic model
288 of divergence via genetic drift, but with a limited number K of mixing events in the history.
289 Our application of Treemix to *Salmonella* gave results which varied in important details
290 depending on the value of K . Each of the events that were identified at a given value of K had
291 counterparts in the inference performed for higher values, but details of the inferred
292 phylogenetic tree and the location and direction of the hybridization events were not
293 consistent. For example, for $K=1$ and $K=2$ Houtenae A and Houtenae B are sister taxa whose
294 common ancestor received genetic material from VII, while for $K=3$, VII and Houtenae B
295 share a common ancestor, which contributed genetic material to Houtenae A. We present the
296 Treemix results for $K=3$ (**Figure 5**) because all of the events inferred are supported by signals
297 identified by chromosome painting and cumulative divergence (**Figures 2, 3 and 4**). The
298 Treemix results with $K=3$ imply that Houtenae A, Houtenae B and VII all have hybrid
299 origins. All three of them received DNA from a shared lineage which branched between
300 Arizonae and Diarizonae (black arrowhead, **Figure 5**), but differ in the remaining source of
301 their ancestry (red arrows, **Figure 5**), which, according to the Treemix estimates, account for
302 about half of their genome in all three cases (1: ancestor of Arizonae to VII: 0.461; 2: Enterica
303 B to Houtenae B: 0.42; 3: ancestor of VII to Houtenae A: 0.49). Note that according to this
304 reconstruction, no pure, or nearly pure, representative of this shared ancestral lineage is
305 present in the sample, a feature which is likely to have contributed largely to the instability of
306 the Treemix analysis and makes all types of evolutionary reconstruction considerably more
307 challenging.

308 The second source for Houtenae B is inferred to be Enterica B (red arrow 2, **Figure 5**), which
309 is consistent with the results from chromosome painting and of the pairwise distances, as
310 discussed above, and is consistent with recent genetic exchange having taken place. The
311 second source for VII is inferred to branch at the same point as Arizonae does in the tree. The
312 deep position of this ancestry source is supported by the distribution of pairwise distances VII
313 has to shallower branching lineages such as Diarizonae or Salamae A, which are more similar

2019-02-14

314 to the distribution found for Arizonae than to that of either Houtenae A or Houtenae B (*e.g.*,
315 **Figure 3C**). The distribution of distances to Arizonae is similar to that of other shallow-
316 branching lineages, suggesting that the recombination was not with Arizonae itself. Finally,
317 the second source for Houtenae A branches next to Salamae A, which is consistent with the
318 reconstructed position of Houtenae A in the phylogenetic tree in Figure 1 and the high
319 coancestry of Houtenae A and Salamae A in Figure 4. However, unlike for Houtenae B, there
320 is no signal of recent recombination of Houtenae A with other lineages in Figure 2.
321 Furthermore, the pairwise distance curves of Salamae A to Houtenae A and Houtenae B are
322 comparable (**Figure 3C**). These features imply that there has not been recent recombination
323 between Houtenae A and Salamae A. Instead, they are consistent with the second source that
324 contributed to Houtenae A being an unsampled sister taxa to Salamae A.

325

326 **Unequal evolutionary rates of the different taxa**

327 One important feature of the phylogenetic tree (**Figure 1**) is the different branch lengths
328 leading to each phylogroup. This feature might be caused by either unequal substitution rates
329 between lineages or by recombination, which can cause hybrid lineages to branch closer to the
330 root. Evidence for unequal substitution rates comes for example from comparisons with
331 Bongori or Arizonae, which can tentatively be treated as outgroups. Salamae A and Salamae
332 B have smaller genetic distances than other lineages to either (**Figure 3D**), despite the
333 chromosome painting indicating no evidence of elevated recombination between them.
334 Furthermore, Salamae A and Salamae B show low genetic distances compared to potential
335 sister lineages to all taxa, suggesting that they have substantially lower substitution rates than
336 other groups. Because our reconstruction of *Salmonella* evolutionary history is incomplete
337 and uncertain, we do not attempt to formally model all of these processes occurring together.

338

339 **Accessory genome relationships**

340 Accessory genes contribute most to ecological specialization and the pattern of horizontal
341 gene transfer among phylogroups might provide important complementary information
342 regarding functional and ecological correlates of the recombination history that we inferred in
343 this work [9]. We therefore analysed the pan-genome of the dataset, which with a protein
344 identity cut-off of 85% rendered a core genome of 1818 gene clusters and a total pan-genome
345 of 21973 genes. Unfortunately, estimations of the strain relationships based on gene

2019-02-14

346 presence/absence and analysis of the shared ancestry revealed that the analyses were strongly
347 affected by the fragmentation of the genomic assemblies (**Table S2**), as was particularly
348 visible for the highly fragmented *Diarizonae* genomes (**Figure S3**). Analysis of the horizontal
349 gene transfer pattern among phylogroups therefore requires higher quality assemblies and will
350 be the subject of future studies.

351

352 **Conclusions**

353 We investigated the diversification and hybridization history within *Salmonella*, a group of
354 prominent public health importance and an early model for microbial speciation and
355 evolutionary studies. By sampling largely in the non-enterica subspecies, we uncovered three
356 novel phylogenetic groups that had not been recognized since the last group, VII, was
357 described in 1991. Our snapshot of diversity within phylogroups of *Salmonella* implies that
358 recombination among phylogroups is relatively rare at any point in time but that when it
359 happens it can be with distantly related lineages rather than sister taxa and can involve large
360 fractions of the core genome. These events are likely to provide substantial potential for
361 phenotypic innovation but may also entail a great deal of hybrid disgenesis.

362

363 The three hybridization events that we have been able to elucidate with any degree of
364 certainty are ongoing or took place in the recent past and all involved a lineage that is not
365 present in unhybridized form in the dataset. This circumstance makes it challenging to
366 estimate simple properties of the events such as the direction of hybridization and the
367 proportion of genome acquired from each source. We can nevertheless robustly conclude that
368 the hybridization has involved at least three entirely different branches of the *Salmonella* tree
369 and has led to the formation of three phylogroups, namely Houtenae A, Houtenae B and VII.
370 Interestingly the latter group was inferred to be a ‘hybrid’ in early MLEE studies [24]. This
371 suggests an interesting question that is likely to be informative about the general nature of
372 species boundaries in bacteria, namely what has happened to make one lineage particularly
373 prone to hybridization in the recent past?

374

375 We see less conclusive but nevertheless still strong evidence for hybridization events in the
376 more distant past. Phylogenetic trees of *Salmonella* phylogroups are notoriously unstable,
377 including in different analyses we have performed (data not shown). In particular,

2019-02-14

378 relationships amongst Salamae A, Salamae B, Diarizonae, Enterica and VIII are difficult to
379 elucidate. The coancestry relationships between these lineages are highly intransitive (**Figure**
380 **4**). One possibility is that this intransitivity is due to a complex multi-way speciation event
381 [5], such that there is no true splitting order to infer. However, it may also represent
382 hybridization events after stepwise speciation. The two lineages that branch deeply (**Figure**
383 **1**), namely VIII and Diarizonae, both show evidence of shared ancestry with basal lineages,
384 Bongori and Arizonae, respectively (**Figure 4**), which is likely to have substantially affected
385 their phylogenetic position.

386

387 The events of recombination inferred in this work explain the difficulties to reconstruct the
388 phylogeny of the genus that have led to multiple distinct hypotheses on the phylogenetic
389 relationships among subspecies. The phylogenetic relationships which do appear to be
390 reasonably certain are that Bongori split from the other phylogroups first, followed by
391 Arizonae and that Indica is a sister group to Enterica. Houtenae A seems to have been a sister
392 taxon of Salamae A, prior to its mixture event. These examples demonstrate that in the right
393 circumstances, phylogenetic signal can be preserved over long evolutionary time periods
394 despite recombination between phylogroups. The problem of reconstructing ancestral
395 hybridization events is a hard one and we do not have the tools or genomes available to
396 reconstruct an entire history with any degree of confidence.

397

398 Our results demonstrate that bacterial species histories are complex. There is considerable
399 phylogenetic signal in the data, consistent with the evolution and long-term persistence
400 barriers to gene flow between lineages but also examples for hybridization events that may
401 reverse species boundaries, sometimes between taxa separated by large genetic distances,
402 rather than between sister taxa. These results mean that phylogenetic trees displaying
403 relationships between species will often represent considerable simplifications of evolutionary
404 history and in the worst case can be entirely misleading. Further work in multiple taxa will
405 elucidate the evolutionary and ecological factors that precipitate speciation and hybridization
406 events.

407

2019-02-14

408 **Author contributions**

409 *Conceptualization: SB, DF. Supervision: SB, DF, NRT, FXW. Performed the experiments:*
410 *SIJ. Data curation: JH, SB, SIJ, FXW. Data analysis: AC, SB, XD, KT. Writing – original*
411 *draft: AC, SB, DF. Writing – review & editing: all.*

412

413 **Conflicts of interest**

414 The author(s) declare that there are no conflicts of interest.

415

416 **Funding information**

417 This work was supported financially by a grant from Region Ile De France to AC and by a
418 Walton Visiting Scientist grant from the Science Foundation of Ireland to SB.

419

420 **Acknowledgements**

421 We acknowledge the expert help of Laure Diancourt, Coralie Tran and Virginie Passet
422 (Institut Pasteur) for MLST data production, and of Mark Achtman for his support at the start
423 of the project and for bioinformatics assistance in MLST data curation.

424

425 **References**

- 426 1. Hugenholtz P, Skarshewski A, Parks DH. Genome-Based Microbial Taxonomy
427 Coming of Age. *Cold Spring Harb Perspect Biol.* **2016**; 8(6).
- 428 2. Adam PS, Borrel G, Brochier-Armanet C, Gribaldo S. The growing tree of
429 Archaea: new perspectives on their diversity, evolution and ecology. *ISME J.* **2017**;
430 11(11):2407–2425.
- 431 3. Brown EW, Mammel MK, LeClerc JE, Cebula TA. Limited boundaries for
432 extensive horizontal gene transfer among Salmonella pathogens. *Proc Natl Acad Sci USA.*
433 **2003**; 100(26):15676–15681.
- 434 4. Octavia S, Lan R. Frequent recombination and low level of clonality within
435 Salmonella enterica subspecies I. *Microbiology (Reading, Engl).* **2006**; 152(Pt 4):1099–1108.
- 436 5. Falush D, Torpdahl M, Didelot X, Conrad DF, Wilson DJ, Achtman M. Mismatch
437 induced speciation in Salmonella: model and data. *Philosophical transactions of the Royal*
438 *Society of London.* **2006**; 361(1475):2045–53.
- 439 6. Sangal V, Harbottle H, Mazzoni CJ, et al. Evolution and population structure of
440 Salmonella enterica serovar Newport. *Journal of bacteriology.* **2010**; 192(24):6465–76.
- 441 7. Didelot X, Bowden R, Street T, et al. Recombination and population structure in
442 Salmonella enterica. *PLoS Genet.* **2011**; 7(7):e1002191.

2019-02-14

- 443 8. Achtman M, Wain J, Weill F-X, et al. Multilocus Sequence Typing as a
444 Replacement for Serotyping in *Salmonella enterica*. *PLoS Pathog.* **2012**; 8(6):e1002776.
- 445 9. Desai PT, Porwollik S, Long F, et al. Evolutionary Genomics of *Salmonella*
446 *enterica* Subspecies. *MBio.* **2013**; 4(2).
- 447 10. Zahrt TC, Maloy S. Barriers to recombination between closely related bacteria:
448 MutS and RecBCD inhibit recombination between *Salmonella typhimurium* and *Salmonella*
449 *typhi*. *Proc Natl Acad Sci USA.* **1997**; 94(18):9786–9791.
- 450 11. Vulić M, Dionisio F, Taddei F, Radman M. Molecular keys to speciation: DNA
451 polymorphism and the control of genetic exchange in enterobacteria. *Proc Natl Acad Sci*
452 *USA.* **1997**; 94(18):9763–9767.
- 453 12. Hanage WP, Fraser C, Spratt BG. The impact of homologous recombination on the
454 generation of diversity in bacteria. *Journal of theoretical biology.* **2006**; 239(2):210–9.
- 455 13. Chen L, Mathema B, Pitout JDD, DeLeo FR, Kreiswirth BN. Epidemic *Klebsiella*
456 *pneumoniae* ST258 is a hybrid strain. *MBio.* **2014**; 5(3):e01355-01314.
- 457 14. Didelot X, Achtman M, Parkhill J, Thomson NR, Falush D. A bimodal pattern of
458 relatedness between the *Salmonella Paratyphi A* and *Typhi* genomes: convergence or
459 divergence by homologous recombination? *Genome research.* **2007**; 17(1):61–8.
- 460 15. Sheppard SK, McCarthy ND, Falush D, Maiden MCJ. Convergence of
461 *Campylobacter* species: implications for bacterial evolution. *Science.* **2008**; 320(5873):237–
462 239.
- 463 16. Brenner FW, Villar RG, Angulo FJ, Tauxe R, Swaminathan B. *Salmonella*
464 nomenclature. *J Clin Microbiol.* **2000**; 38(7):2465–2467.
- 465 17. Tindall BJ, Grimont PA, Garrity GM, Euzéby JP. Nomenclature and taxonomy of
466 the genus *Salmonella*. *Int J Syst Evol Microbiol.* **2005**; 55(Pt 1):521–4.
- 467 18. Reeves MW, Evins GM, Heiba AA, Plikaytis BD, Farmer JJ. Clonal nature of
468 *Salmonella typhi* and its genetic relatedness to other salmonellae as shown by multilocus
469 enzyme electrophoresis, and proposal of *Salmonella bongori* comb. nov. *J Clin Microbiol.*
470 **1989**; 27(2):313–320.
- 471 19. Parry CM, Hien TT, Dougan G, White NJ, Farrar JJ. Typhoid fever. *N Engl J Med.*
472 **2002**; 347(22):1770–1782.
- 473 20. Sánchez-Vargas FM, Abu-El-Haija MA, Gómez-Duarte OG. *Salmonella*
474 infections: an update on epidemiology, management, and prevention. *Travel Med Infect Dis.*
475 **2011**; 9(6):263–277.
- 476 21. Lamas A, Miranda JM, Regal P, Vázquez B, Franco CM, Cepeda A. A
477 comprehensive review of non-enterica subspecies of *Salmonella enterica*. *Microbiol Res.*
478 **2018**; 206:60–73.
- 479 22. Selander R, Beltran P, Smith N. Evolutionary genetics of *Salmonella*. Evolution at
480 the molecular level. Selander RK, Clark AG, Whittam TS, editors; 1991. p. 25–57.
- 481 23. Nelson K, Selander RK. Evolutionary genetics of the proline permease gene (*putP*)
482 and the control region of the proline utilization operon in populations of *Salmonella* and
483 *Escherichia coli*. *J Bacteriol.* **1992**; 174(21):6886–6895.
- 484 24. Boyd EF, Wang FS, Whittam TS, Selander RK. Molecular genetic relationships of
485 the salmonellae. *Applied and environmental microbiology.* **1996**; 62(3):804–8.
- 486 25. Guibourdenche M, Roggentin P, Mikoleit M, et al. Supplement 2003-2007 (No. 47)
487 to the White-Kauffmann-Le Minor scheme. *Research in microbiology.* **2010**; 161(1):26–9.
- 488 26. Shelobolina ES, Sullivan SA, O'Neill KR, Nevin KP, Lovley DR. Isolation,
489 characterization, and U(VI)-reducing potential of a facultatively anaerobic, acid-resistant
490 bacterium from Low-pH, nitrate- and U(VI)-contaminated subsurface sediment and
491 description of *Salmonella subterranea* sp. nov. *Appl Environ Microbiol.* **2004**; 70(5):2959–
492 2965.

2019-02-14

- 493 27. Nelson K, Whittam TS, Selander RK. Nucleotide polymorphism and evolution in
494 the glyceraldehyde-3-phosphate dehydrogenase gene (*gapA*) in natural populations of
495 *Salmonella* and *Escherichia coli*. *Proc Natl Acad Sci U S A*. **1991**; 88(15):6667–71.
- 496 28. Thampapillai G, Lan R, Reeves PR. Molecular evolution in the *gnd* locus of
497 *Salmonella enterica*. *Mol Biol Evol*. **1994**; 11(6):813–28.
- 498 29. Christensen H, Nordentoft S, Olsen JE. Phylogenetic relationships of *Salmonella*
499 based on rRNA sequences. *Int J Syst Bacteriol*. **1998**; 48 Pt 2:605–610.
- 500 30. Brown EW, Kotewicz ML, Cebula TA. Detection of recombination among
501 *Salmonella enterica* strains using the incongruence length difference test. *Mol Phylogenet*
502 *Evol*. **2002**; 24(1):102–120.
- 503 31. Whittam TS, Bumbaugh AC. Inferences from whole-genome sequences of bacterial
504 pathogens. *Curr Opin Genet Dev*. **2002**; 12(6):719–725.
- 505 32. Porwollik S, Wong RM-Y, McClelland M. Evolutionary genomics of *Salmonella*:
506 gene acquisitions revealed by microarray analysis. *Proc Natl Acad Sci USA*. **2002**;
507 99(13):8956–8961.
- 508 33. McQuiston JR, Herrera-Leon S, Wertheim BC, et al. Molecular phylogeny of the
509 salmonellae: relationships among *Salmonella* species and subspecies determined from four
510 housekeeping genes and evidence of lateral gene transfer events. *Journal of bacteriology*.
511 **2008**; 190(21):7060–7.
- 512 34. Fookes M, Schroeder GN, Langridge GC, et al. *Salmonella bongori* provides
513 insights into the evolution of the Salmonellae. *PLoS Pathog*. **2011**; 7(8):e1002191.
- 514 35. Trujillo S, Keys CE, Brown EW. Evaluation of the taxonomic utility of six-enzyme
515 pulsed-field gel electrophoresis in reconstructing *Salmonella* subspecies phylogeny. *Infect*
516 *Genet Evol*. **2011**; 11(1):92–102.
- 517 36. Pettengill JB, Timme RE, Barrangou R, et al. The evolutionary history and
518 diagnostic utility of the CRISPR-Cas system within *Salmonella enterica* ssp. *enterica*. *PeerJ*.
519 **2014**; 2:e340.
- 520 37. Kisiela DI, Chattopadhyay S, Libby SJ, et al. Evolution of *Salmonella enterica*
521 virulence via point mutations in the fimbrial adhesin. *PLoS Pathog*. **2012**; 8(6):e1002733.
- 522 38. Kidgell C, Reichard U, Wain J, et al. *Salmonella typhi*, the causative agent of
523 typhoid fever, is approximately 50,000 years old. *Infect Genet Evol*. **2002**; 2(1):39–45.
- 524 39. Criscuolo A, Gascuel O. Fast NJ-like algorithms to deal with incomplete distance
525 matrices. *BMC Bioinformatics*. **2008**; 9:166.
- 526 40. Nei M, Li WH. Mathematical model for studying genetic variation in terms of
527 restriction endonucleases. *Proc Natl Acad Sci USA*. **1979**; 76(10):5269–5273.
- 528 41. Librado P, Rozas J. DnaSP v5: a software for comprehensive analysis of DNA
529 polymorphism data. *Bioinformatics*. **2009**; 25(11):1451–2.
- 530 42. McClelland M, Sanderson KE, Spieth J, et al. Complete genome sequence of
531 *Salmonella enterica* serovar Typhimurium LT2. *Nature*. **2001**; 413(6858):852–6.
- 532 43. Altschul SF, Madden TL, Schaffer AA, et al. Gapped BLAST and PSI-BLAST: a
533 new generation of protein database search programs. *Nucleic acids research*. **1997**;
534 25(17):3389–402.
- 535 44. Lefort V, Desper R, Gascuel O. FastME 2.0: A Comprehensive, Accurate, and Fast
536 Distance-Based Phylogeny Inference Program. *Mol Biol Evol*. **2015**; 32(10):2798–2800.
- 537 45. Lawson DJ, Hellenthal G, Myers S, Falush D. Inference of population structure
538 using dense haplotype data. *PLoS Genet*. **2012**; 8(1):e1002453.
- 539 46. Pickrell JK, Pritchard JK. Inference of population splits and mixtures from
540 genome-wide allele frequency data. *PLoS Genet*. **2012**; 8(11):e1002967.
- 541 47. Page AJ, Cummins CA, Hunt M, et al. Roary: rapid large-scale prokaryote pan
542 genome analysis. *Bioinformatics*. **2015**; 31(22):3691–3693.

2019-02-14

- 543 48. Wei, T, Simko V. R package “corrplot”: Visualization of a Correlation Matrix
544 (Version 0.84) [Internet]. 2017. Available from: <https://github.com/taiyun/corrplot>
545 49. Alikhan N-F, Zhou Z, Sergeant MJ, Achtman M. A genomic overview of the
546 population structure of Salmonella. PLoS Genet. **2018**; 14(4):e1007261.
547

548

549 **Figures**

550 **Figure 1. Phylogenetic tree of 73 *Salmonella* strains based on all shared core genes.**

551 The balanced minimum-evolution phylogenetic tree was constructed using FastME (see
552 Methods). The 11 phylogroups are indicated above their ancestral branch; Enterica groups A
553 and B are also indicated. Bootstrap/branch support values are indicated at the nodes.

554

555 **Figure 2. Coancestry matrix of 73 *Salmonella* genomes, computed using** 556 **CHROMOPAINTER.**

557

558 **Figure 3. Cumulative curves of gene-by-gene distances between selected pairs of**
559 **genomes. A:** Comparisons with Enterica (group B, serovar Schwartzengrund CVM19633).
560 The arrowhead shows that 20% (0.20, Y-axis) of the genes of an Enterica B strain have less
561 than 1% (0.01, x-axis) divergence to Houtenae B. **B:** Comparisons with Houtenae B
562 (2193/78). The arrowhead shows that 5% of the VII genome and 6% of Houtenae A has less
563 than 0.1% divergence with Houtenae B. **C:** Comparisons with Salamae A (1268/72). **D:**
564 Comparisons with Arizonae (CDC 129-73).

565

566 **Figure 4. Coancestry matrix between 9 unrelated genomes, computed using** 567 **CHROMOPAINTER.**

568

569 **Figure 5. Treemix analysis of 12 genomes representative of phylogroups diversity.**

570 The arrowhead indicates the position of the ancestor contributing to extant Houtenae A,
571 Houtenae B and VII lineages. The red arrows indicate gene fluxes inferred by Treemix.

2019-02-14

572 **Supporting Information**

573

574 **Table S1. Strains studied by MLST.**

575 **Table S2. Genomic sequence data.**

576 **Table S3. Primers and conditions used for MLST gene amplification and sequencing for**
577 **non-Enterica isolates.**

578

579 **Figure S1. BioNJ* tree of 382 *Salmonella* strains based on seven housekeeping gene**
580 **sequences.** The inset shows the average nucleotide diversity of each phylogroup (houtenae
581 comprises Houtenae A and Houtenae B) at the seven MLST genes.

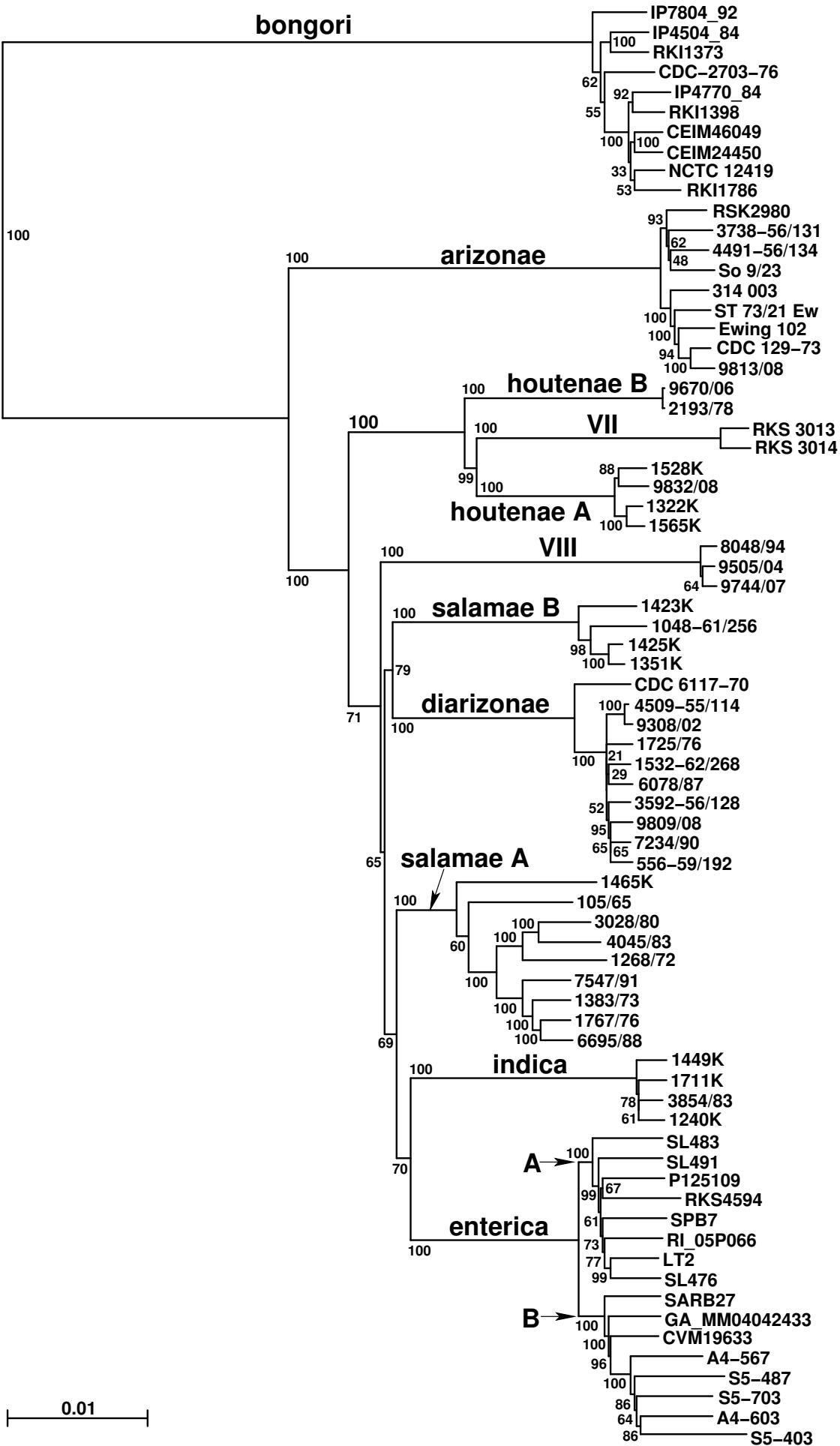
582

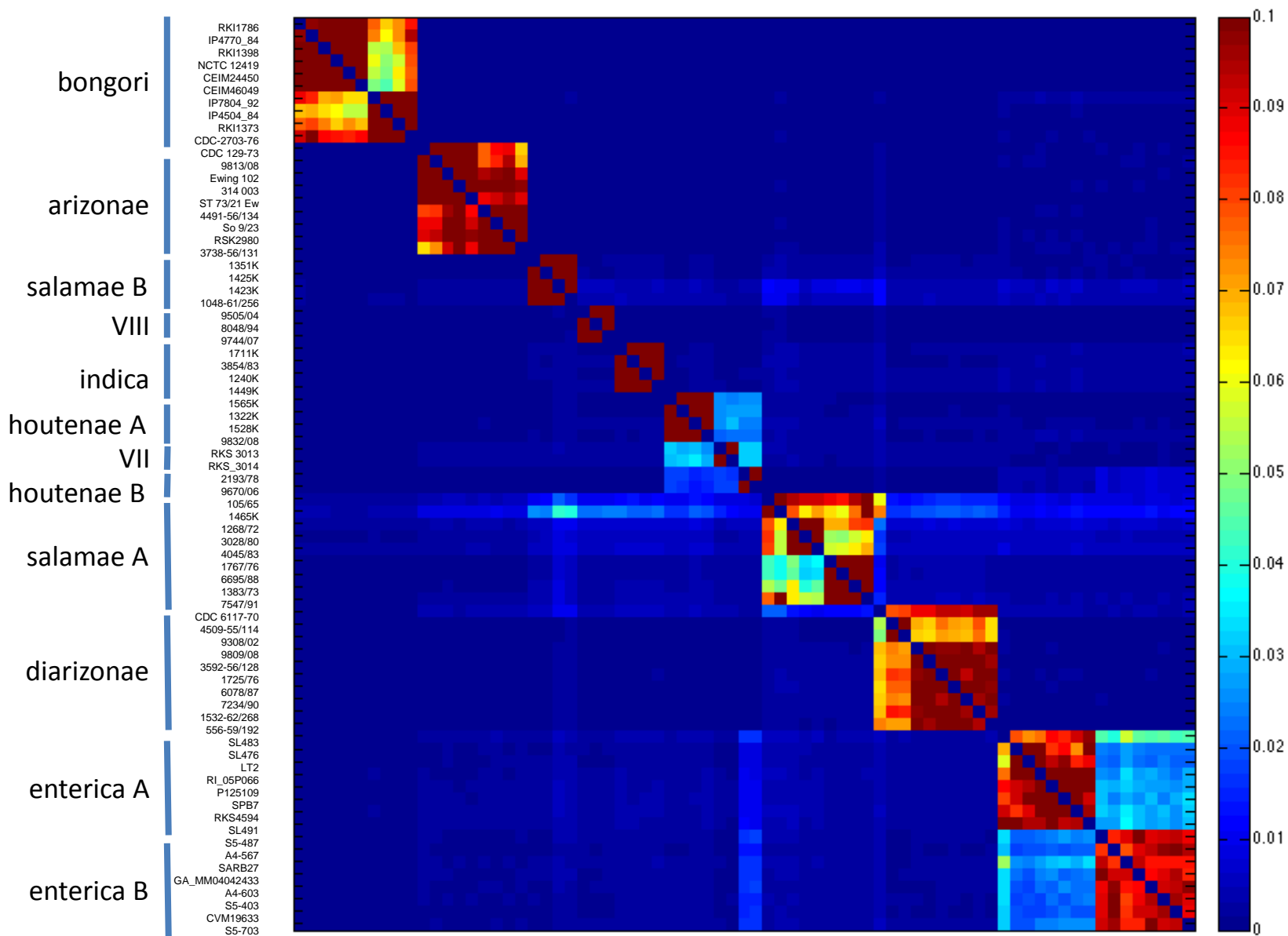
583 **Figure S2. Minimum spanning tree representations of the genotypic diversity within**
584 ***Salmonella* groups.** The minimum spanning trees were constructed for each group based on
585 number of mismatches among MLST allelic profiles. Strains selected for genome sequencing
586 are represented by blue sectors (or blue circles when only one strain shared the corresponding
587 genotype). Grey zones surround groups of sequence types that are connected successively by
588 single allelic mismatches and are equivalent to clonal complexes or ‘eBURST’ groups
589 (Achtman et al., 2012).

590

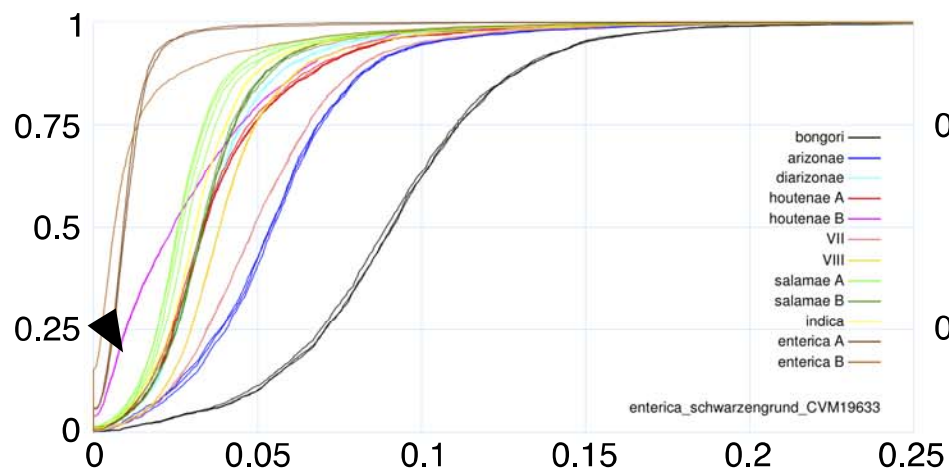
591 **Figure S3. Heatmap of the proportion of shared genes**

592 Strains are ordered according to the phylogeny in Figure 1 (left). The proportion of shared
593 genes was computed from the ROARY output with a protein identity cut-off of 85% and the
594 “don’t split paralogs” option.

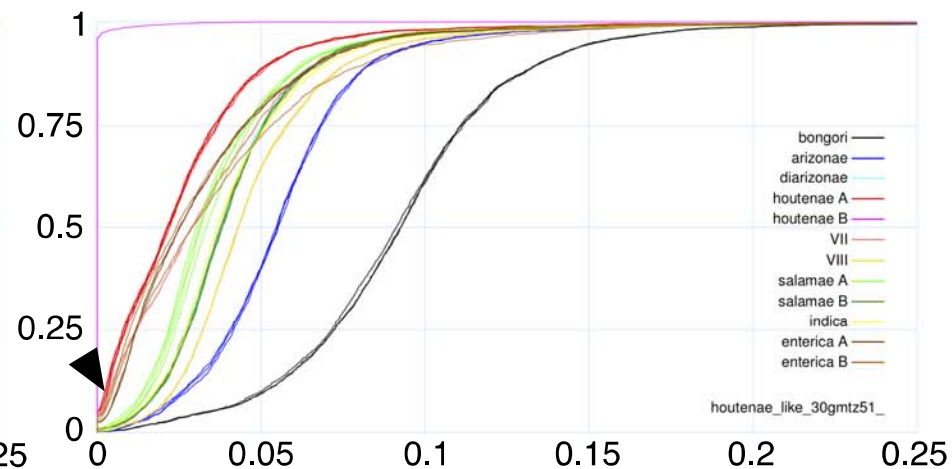




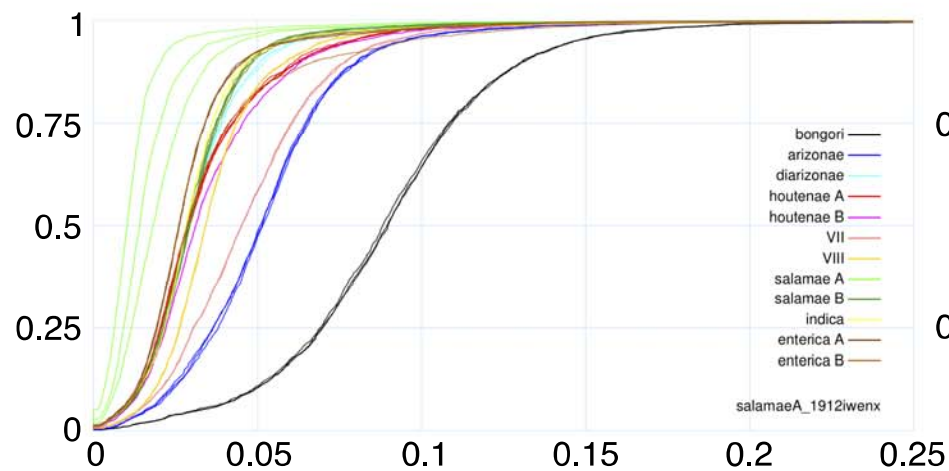
A Enterica



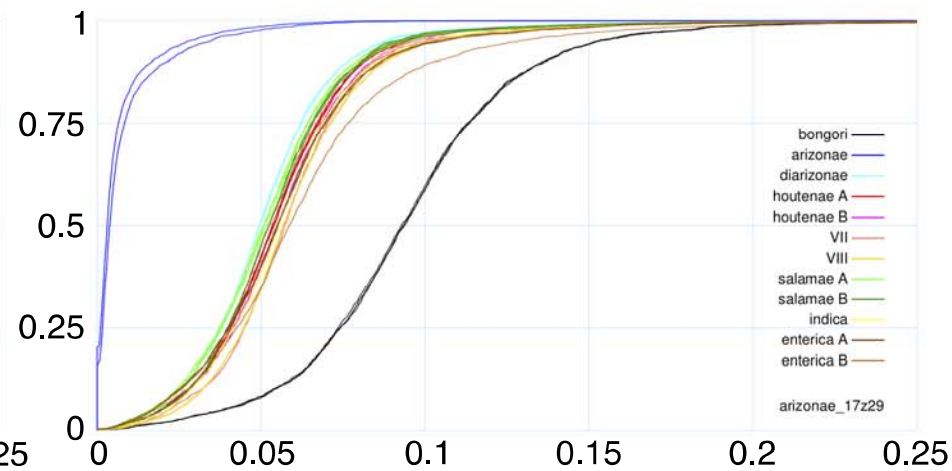
B Houtenae B

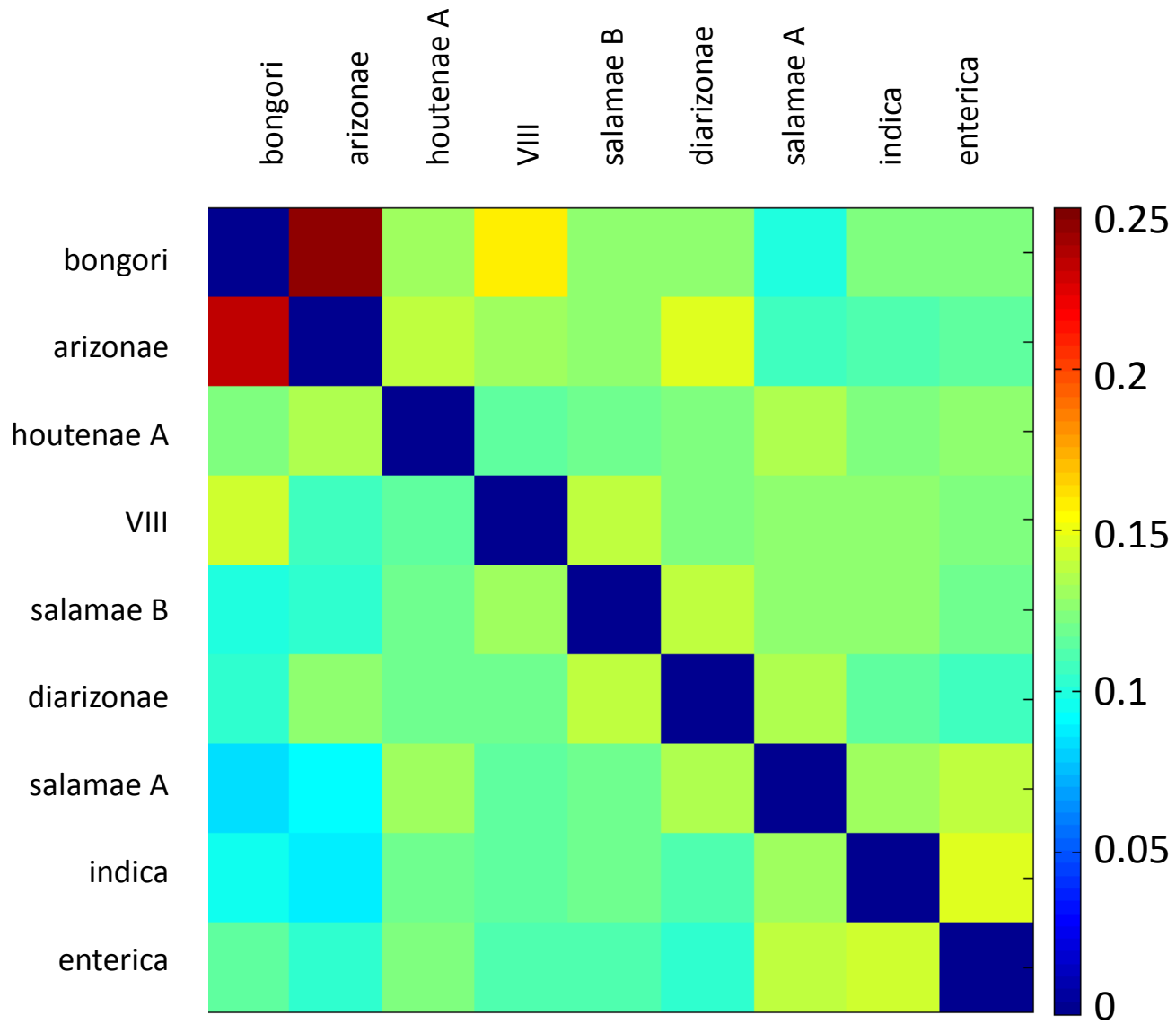


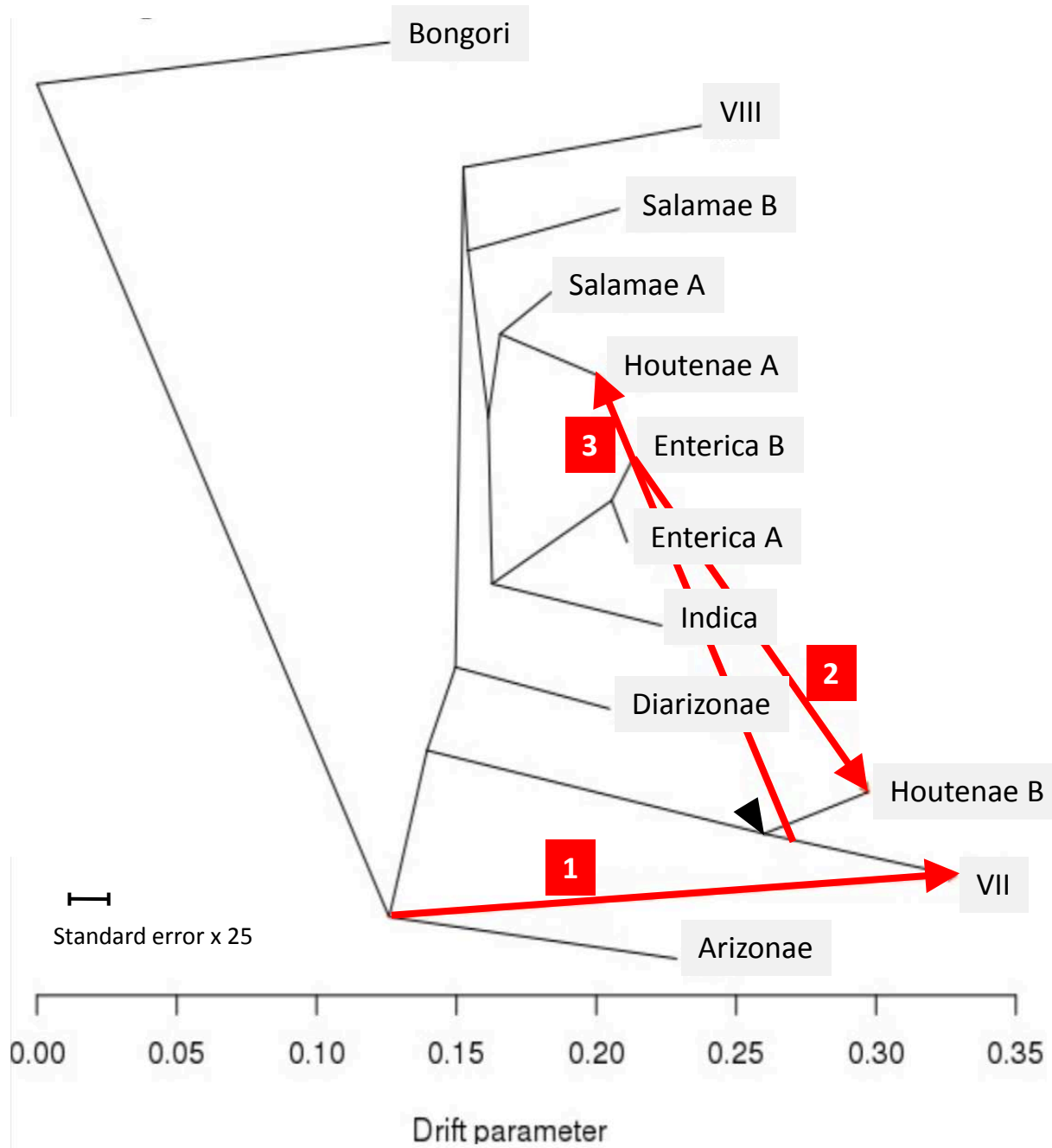
C Salamae A



D Arizonae

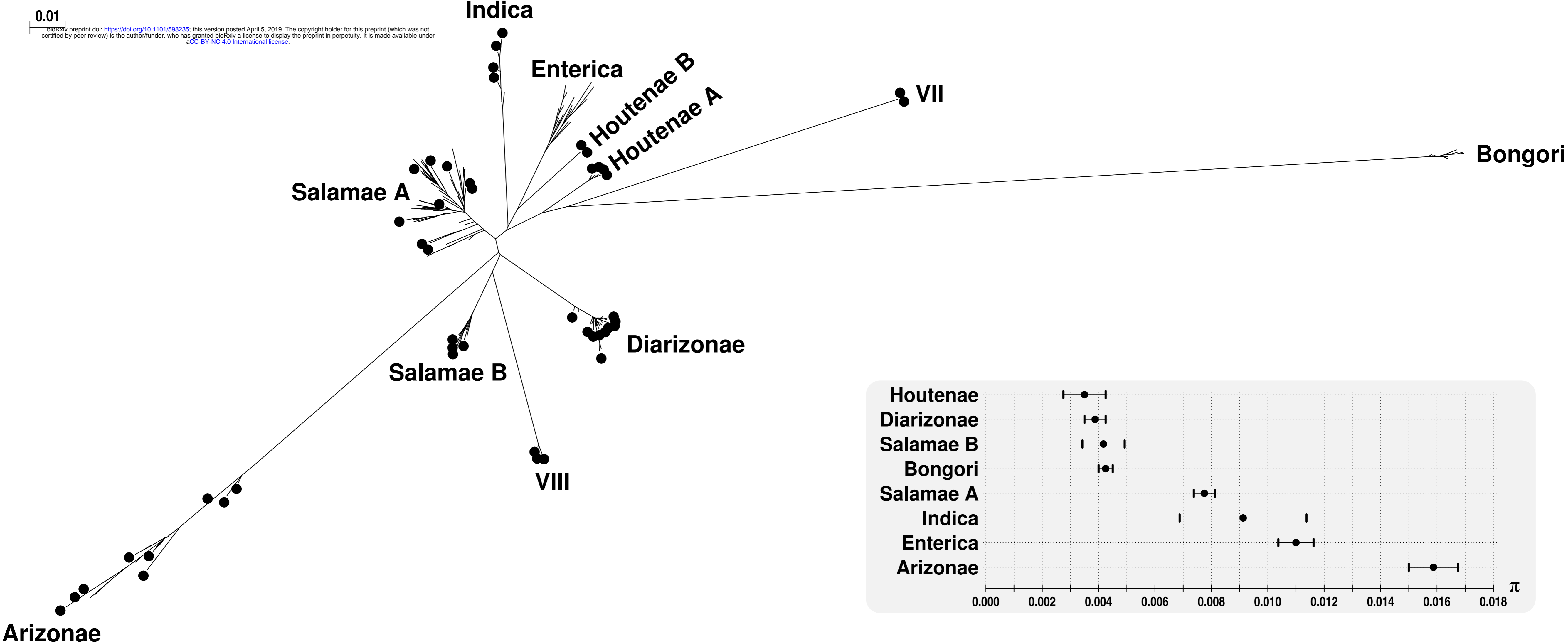




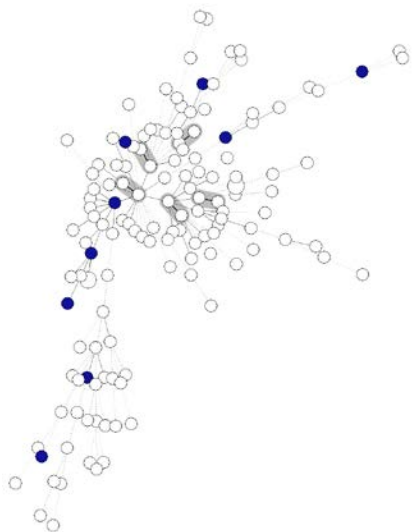


0.01

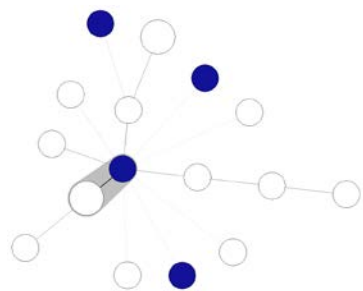
bioRxiv preprint doi: <https://doi.org/10.1101/598235>; this version posted April 5, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



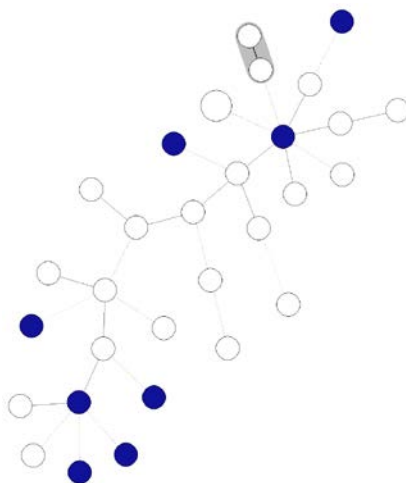
Salamae A



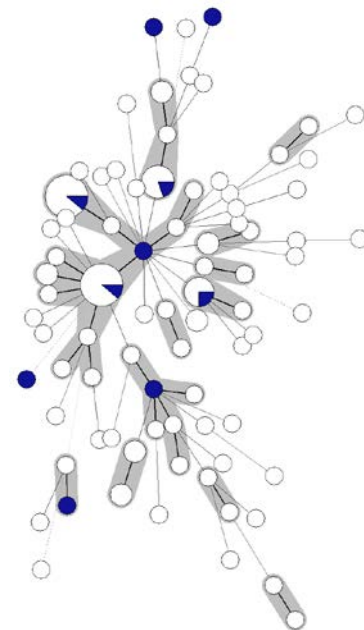
Salamae B



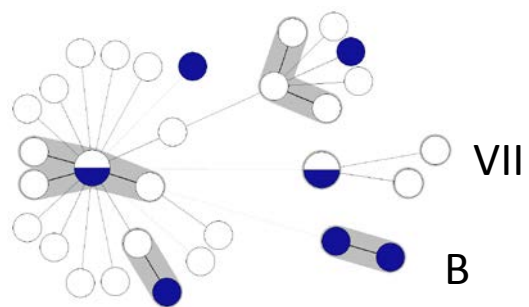
Arizonae



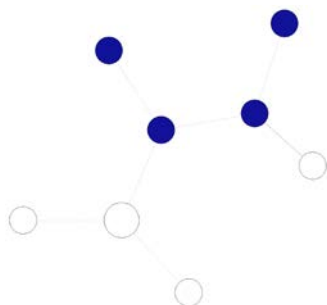
Diarizonae



Houtenae A, B and VII



Indica



VIII



S. bongori



