

1 Supplementary Information:

2 Supplementary Figures 1–16 and Supplementary Table 1

3

4 Title: Apparent nosocomial adaptation of *Enterococcus faecalis* predates the modern
5 hospital era

6

7 Author information:

8 Anna K. Pöntinen^{1*}, Janetta Top², Sergio Arredondo-Alonso^{1,2}, Gerry Tonkin-Hill³, Ana R.

9 Freitas⁴, Carla Novais⁴, Rebecca A. Gladstone¹, Maiju Pesonen⁵, Rodrigo Meneses², Henri

10 Pesonen¹, John A. Lees⁶, Dorota Jamrozy³, Stephen D. Bentley³, Val Fernandez Lanza⁷,

11 Carmen Torres⁸, Luisa Peixe⁴, Teresa M. Coque^{7,9}, Julian Parkhill^{10,11}, Anita C. Schürch², Rob

12 J. L. Willems², Jukka Corander^{1,3,12*}

13

14 * Corresponding authors: Anna K. Pöntinen (a.k.pontinen@medisin.uio.no), Jukka Corander

15 (jukka.corander@medisin.uio.no)

16 These authors contributed equally: Anita C. Schürch, Rob J. L. Willems, Jukka Corander

17

18 Affiliations:

19 ¹ Department of Biostatistics, Faculty of Medicine, University of Oslo, Oslo, Norway

20 ² Department of Medical Microbiology, University Medical Center Utrecht, Utrecht, The
21 Netherlands

22 ³ Parasites and Microbes, Wellcome Sanger Institute, Cambridge, United Kingdom

23 ⁴ UCIBIO/REQUIMTE. Laboratory of Microbiology, Biological Sciences Department, Faculty of
24 Pharmacy, University of Porto, Porto, Portugal

25 ⁵ Oslo Centre for Biostatistics and Epidemiology (OCBE), Oslo University Hospital Research
26 Support Services, Oslo, Norway

27 ⁶ MRC Centre for Global Infectious Disease Analysis, Department of Infectious Disease
28 Epidemiology, Imperial College London, London, United Kingdom

29 ⁷ Department of Microbiology, Ramón y Cajal Institute for Health Research, Ramón y Cajal
30 University Hospital, Madrid, Spain

31 ⁸ Department of Food and Agriculture, Area of Biochemistry and Molecular Biology,
32 University of La Rioja, Logroño, Spain

33 ⁹ CIBER in Epidemiology and Public Health (CIBERESP), Madrid, Spain

34 ¹⁰ Wellcome Sanger Institute, Cambridge, United Kingdom

35 ¹¹ Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom

36 ¹² Helsinki Institute of Information Technology, Department of Mathematics and Statistics,
37 University of Helsinki, Helsinki, Finland

38

39

40

41

42

43

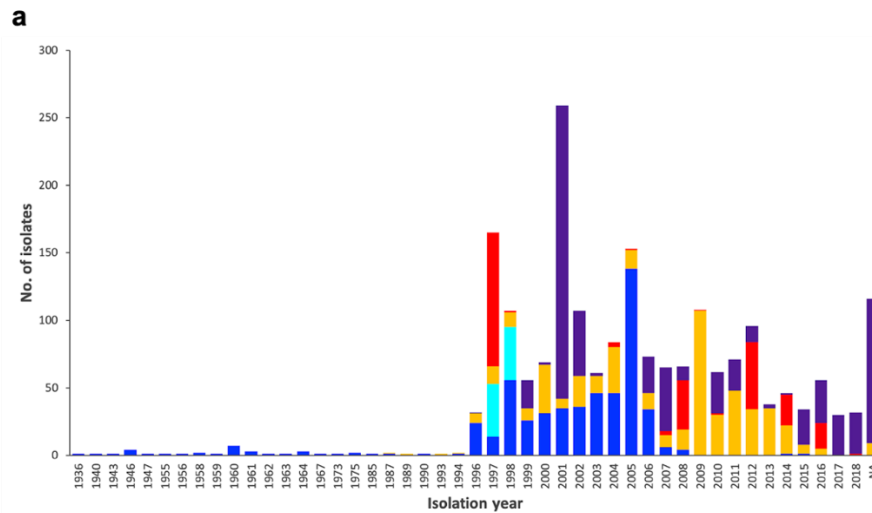
44

45

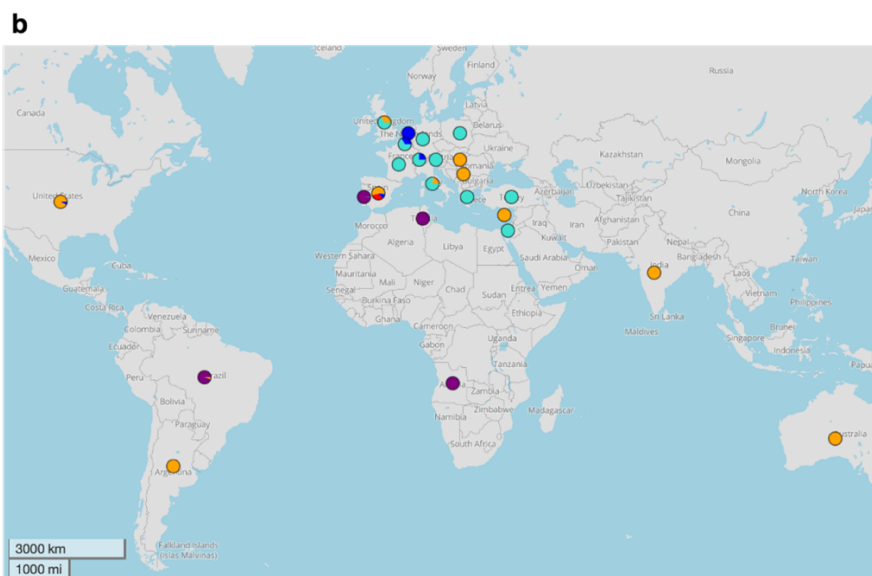
46

47

48



49



50

51 **Supplementary Fig. 1: Temporal and geographic distribution of the *E. faecalis* collection.**

52 Total of 2,027 *E. faecalis* isolates were used in the study, representing collections from

53 University Medical Center Utrecht, The Netherlands ($n = 535$; blue), European Network for

54 Antibiotic Resistance and Epidemiology (ENARE) at the University Medical Center Utrecht,

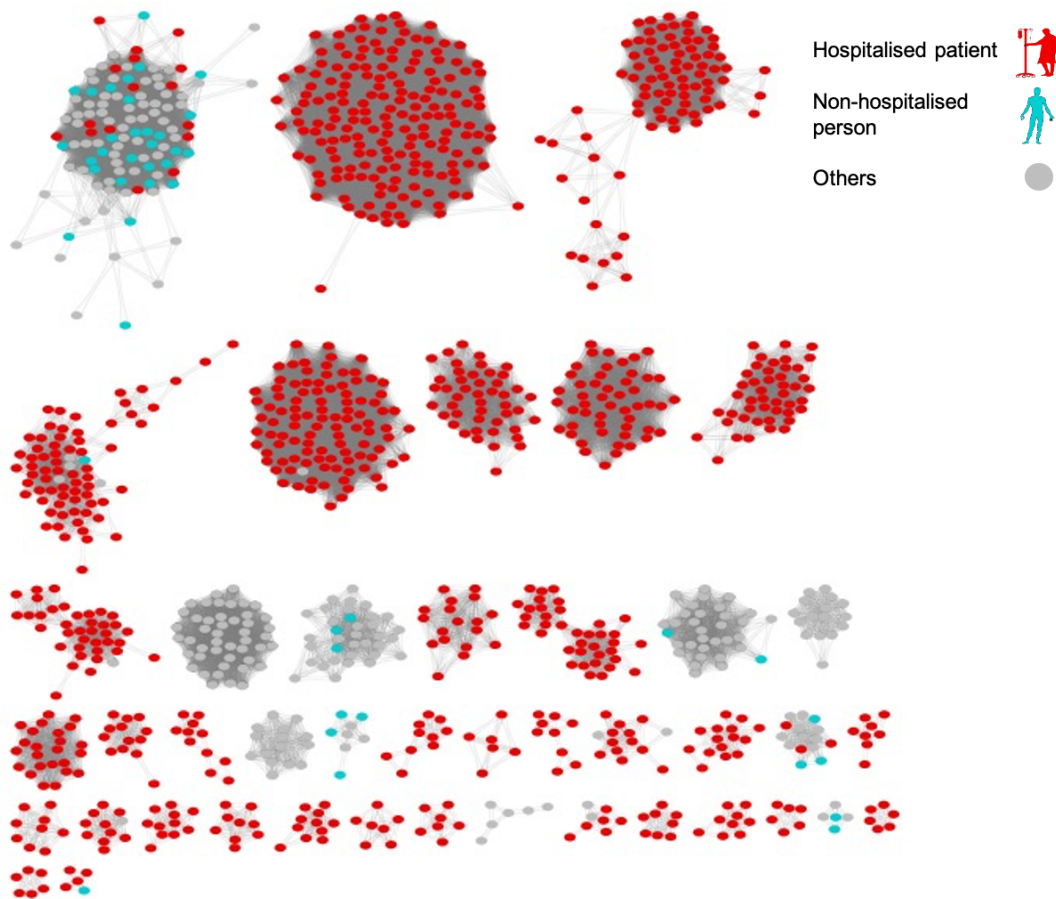
55 The Netherlands ($n = 78$; light blue), University Hospital Ramòn y Cajal, Spain ($n = 503$;

56 orange), University of Porto, Portugal ($n = 671$; purple), and University of La Rioja, Spain ($n =$

57 240; red). **a**, Number of isolates per isolation years and isolate collections. **b**, Geographic

58 distribution of isolates, created by using Microreact¹. Source data are provided as a Source

59 Data file.



60

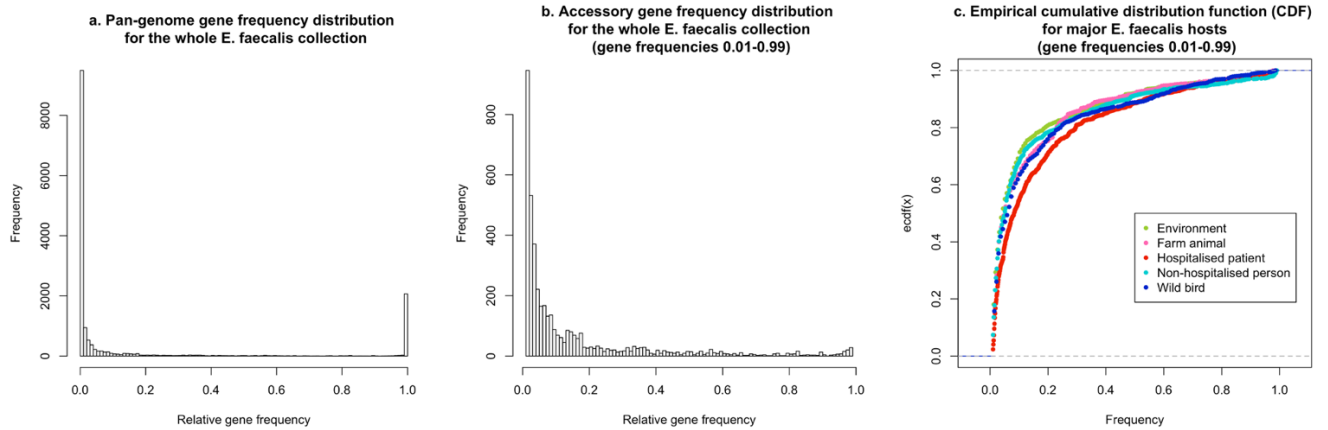
61 **Supplementary Fig. 2: Network analysis of the *E. faecium* accessory genomes, as defined**
 62 **by Panaroo², depicts clearly defined hospital-associated (HA) clusters.** Nodes indicate
 63 separate isolates, connected when shared $\geq 95\%$ of their accessory genome, and colour-
 64 coded according to their isolation sources as indicated in the legend: hospitalised patient
 65 (red), non-hospitalised person (light blue), and others (grey). Components of less than five
 66 isolates were filtered out, and the resulting network was visualised using Cytoscape³. Source
 67 data are provided as a Source Data file.

68

69

70

71



72

73 **Supplementary Fig. 3: Accessory gene frequency distributions present no significant**

74 **differences between different host types, indicating generalist nature for *E. faecalis*. a,**

75 **Pangenome gene frequency histogram for the whole collection of 2,027 *E. faecalis* isolates.**

76 **b,** Accessory gene frequency distribution for the whole collection, with gene frequencies of

77 1–99 %.

78 **c,** Empirical cumulative distribution functions (CDFs) of 1–99% gene frequencies for

79 major *E. faecalis* host types as indicated by colour coding: environment (green), farm animal

80 (pink), hospitalised patient (red), non-hospitalised person (light blue), and wild bird (dark

81 blue). There were no significant differences between empirical CDFs of hospital-associated

82 (HA) isolates and other host types ($P > 0.20$; one-sided permutation tests with the maximum

83 difference of each pair of empirical CDFs as the test statistic).

84

85

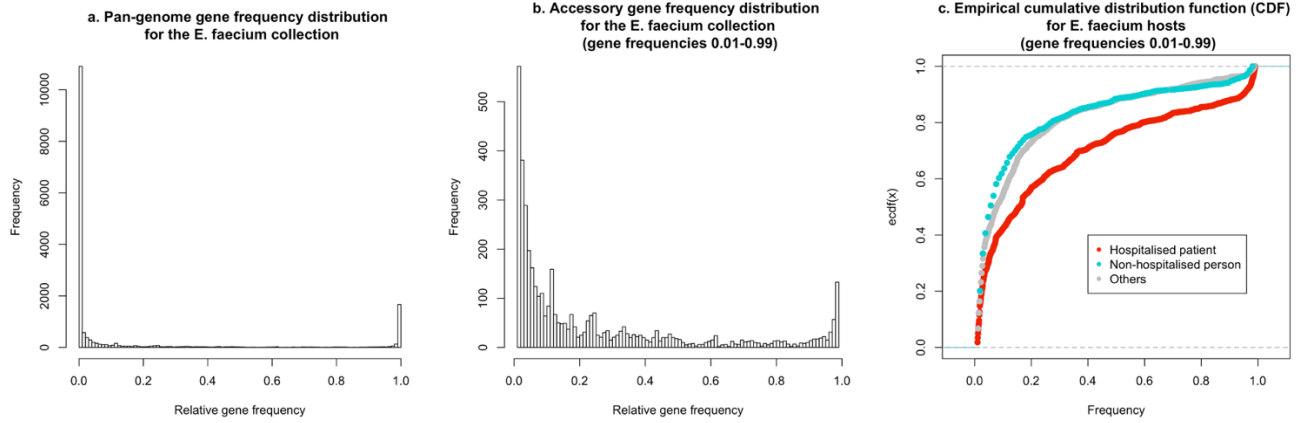
86

87

88

89

90



91

92 **Supplementary Fig. 4: Accessory gene frequency distributions indicate adaptation of**

93 ***E. faecium* isolates of hospital origin. a,** Pangenome gene frequency histogram for the

94 1,602 *E. faecium* isolates⁴. **b,** Accessory gene frequency distribution, with gene frequencies

95 of 1–99 % . **c,** Empirical cumulative distribution functions (CDFs) of 1–99% gene frequencies

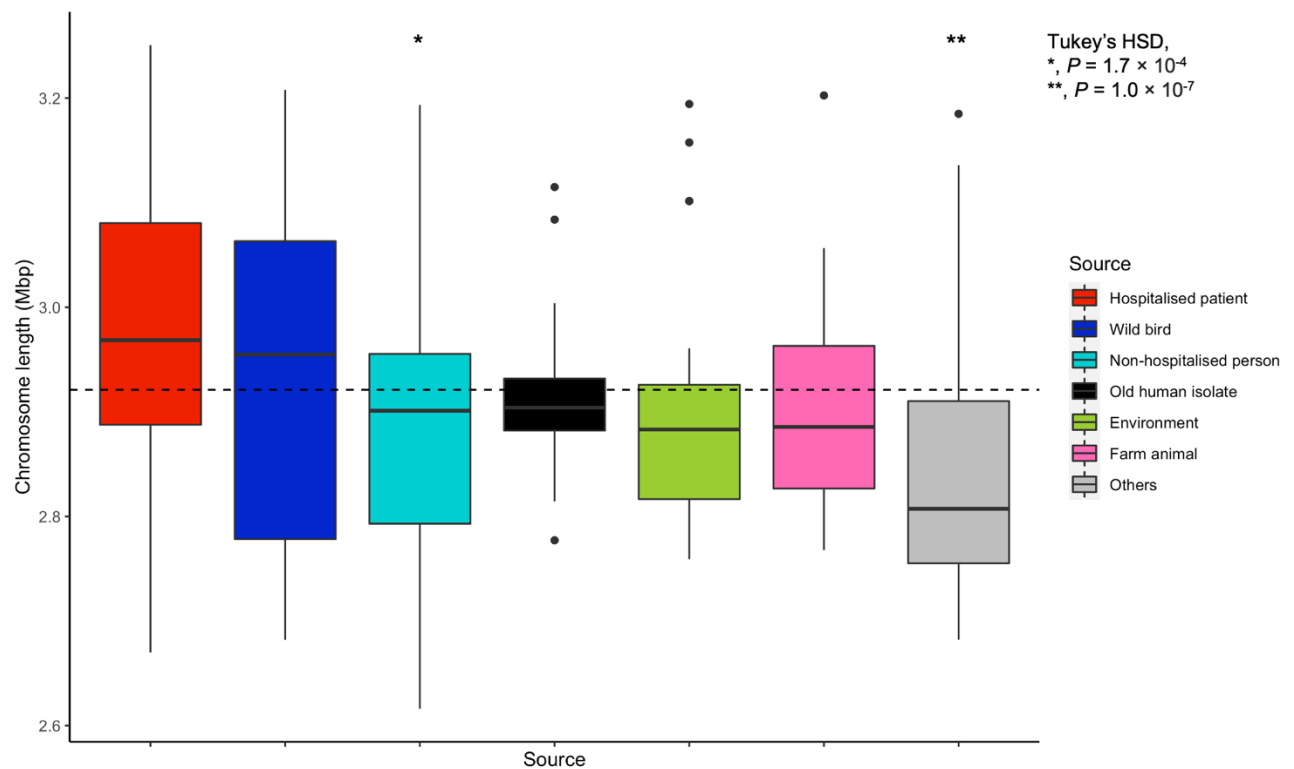
96 for *E. faecium* host types as indicated by colour coding: hospitalised patient (red), non-

97 hospitalised person (light blue), and others (grey). Empirical CDF for hospital-associated (HA)

98 *E. faecium* isolates differed significantly from those of non-hospital origins ($P < 0.05$; one-

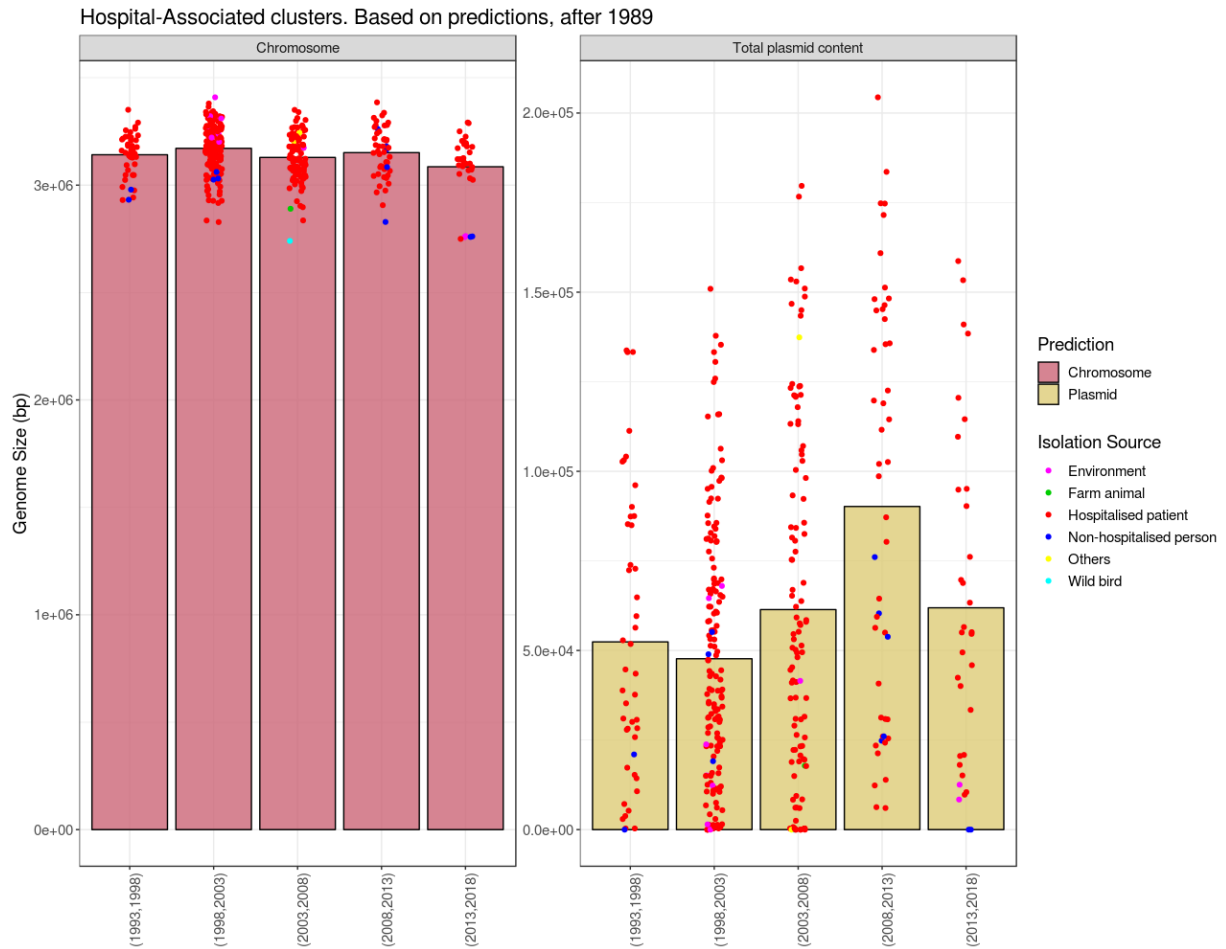
99 sided permutation tests with the maximum difference of each pair of empirical CDFs as the

100 test statistic).



101

102 **Supplementary Fig. 5: Comparisons of hybrid assembly chromosomes show mainly stable**
 103 **chromosome sizes (Mbp) across different host types in *E. faecalis*, reflecting the generalist**
 104 **lifestyle.** Boxplots show the median values of chromosome lengths from isolates colour-
 105 coded as per their isolation source: hospitalised patient (red; $n = 967$), wild bird (dark blue; n
 106 = 136), non-hospitalised person (light blue; $n = 391$), old human isolates (black; $n = 28$),
 107 environment (green; $n = 156$), farm animals (pink; $n = 130$), and others (grey; $n = 219$). First
 108 and third quartiles are shown and whiskers extending at most to the value of $1.5 \times$ inter-
 109 quartile range from the hinge. Individual dots depict single outliers, and the dashed line
 110 indicates mean chromosome length across all sources. *, $P = 1.7 \times 10^{-4}$; **, $P = 1.0 \times 10^{-7}$, as
 111 compared to the mean chromosome size of hospitalised patient isolates; Tukey's HSD,
 112 accounting for multiple comparison. Source data are provided as a Source Data file.



113

114 **Supplementary Fig. 6: Predicted chromosome (pink) sizes (bp) in the hospital-associated**

115 **(HA) clusters show no increase over the years of isolation, while a slight intermittently**

116 **increasing trend is seen in the predicted plasmid content (yellow-green) sizes (bp).**

117 Predictions were derived from mlplasmids⁵. Bar plots represent mean genome sizes (bp),

118 and each node represents a single isolate, coded by colour as indicated in the legend:

119 environment (pink), farm animal (green), hospitalised patient (red), non-hospitalised person

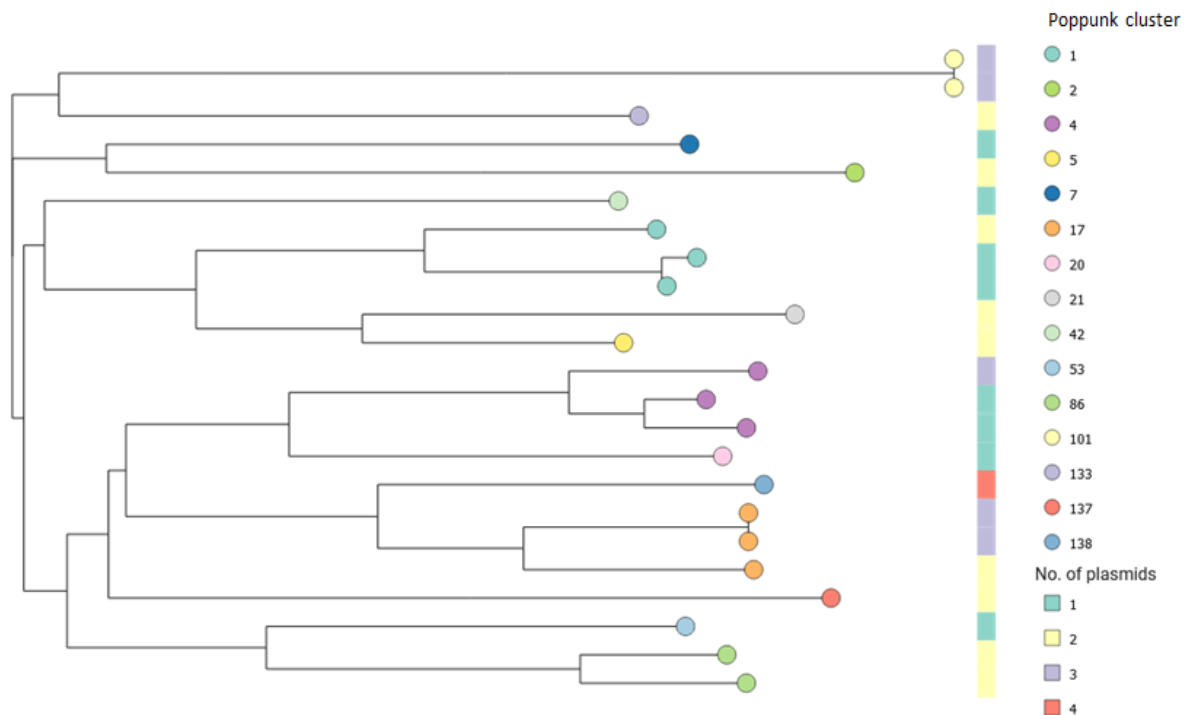
120 (dark blue), others (yellow), and wild bird (light blue). Years are shown in intervals of five

121 years.

122

123

124



125

126 **Supplementary Fig. 7: Alignment-free k-mer based clustering⁶ of 23 old isolates shows**

127 **largely diverse genomic background.** Nodes are labelled by colour, as indicated in the

128 legend, according to their cluster defined by Population Partitioning Using Nucleotide K-

129 mers (PopPUNK)⁷. The number of plasmids for each strain is coded in the dendrogram by

130 colour as indicated in the legend: 1 (turquoise), 2 (yellow), 3 (purple), and 4 (red).

131

132

133

134

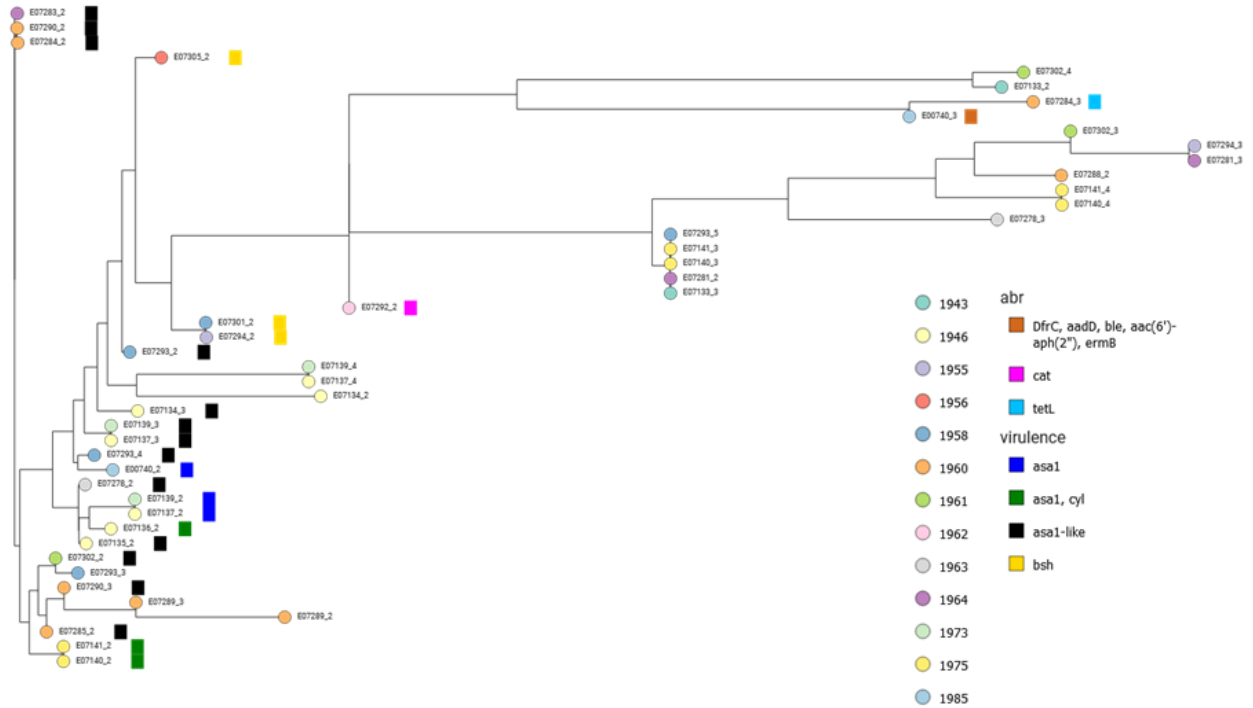
135

136

137

138

139



140

141 **Supplementary Fig. 8: Clustering of presence and absence of k-mers⁶ of all 45 identified**

142 **plasmids among 23 old isolates reveals large plasmid diversity. Nodes are coloured**

143 according to isolation year as indicated in the legend. The presence of antibiotic resistance

144 genes (abr) and/or virulence genes are indicated next to the plasmid name, coloured as

145 indicated in the legend.

146

147

148

149

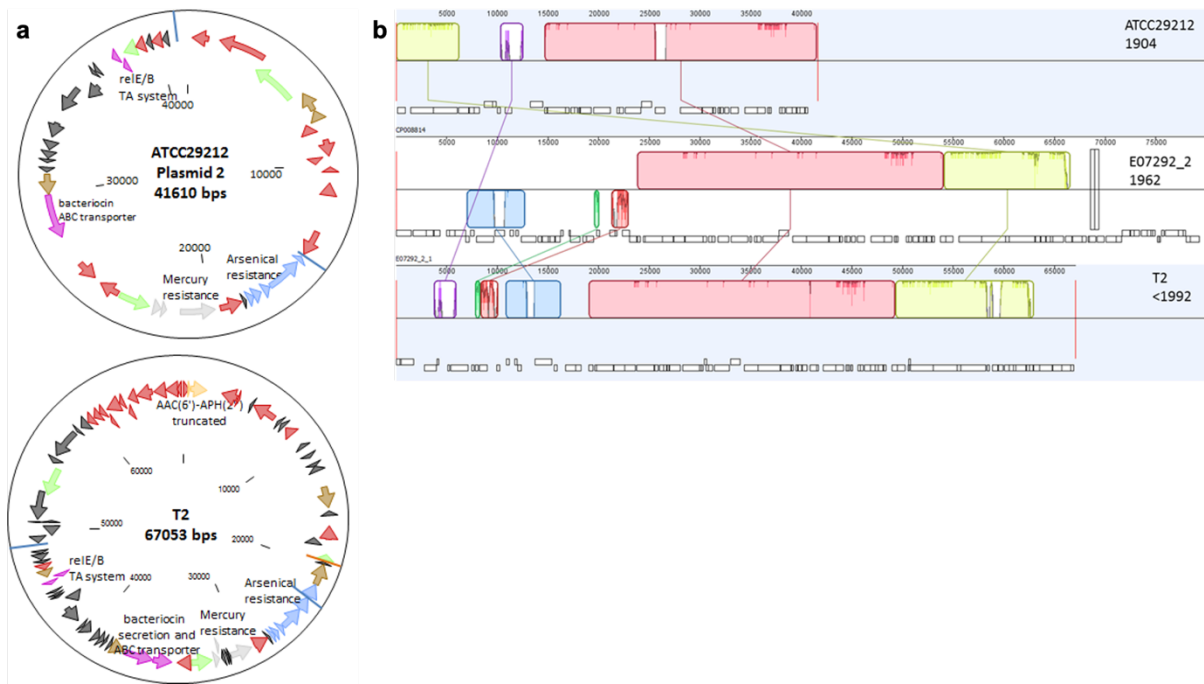
150

151

152

153

154



155

156 **Supplementary Fig. 9: One of the old *E. faecalis* plasmids from 1962 shows identical**

157 **regions to plasmids geographically and temporally widely distributed. a, Genomic**

158 **organization for strain ATCC29212 plasmid 2 and the T2 plasmid, isolated in the UK in 1904**

159 **and in Japan prior to 1992, respectively. Red arrows indicate plasmid associated genes and**

160 **transposases, yellow arrows indicate antimicrobial resistance genes, light grey mercury**

161 **resistance and light blue arsenical resistance genes. b, MAUVE alignment⁸ of the plasmids**

162 **from strain ATCC29212, E07292, and T2, indicating identical regions shown in red and**

163 **yellow.**

164

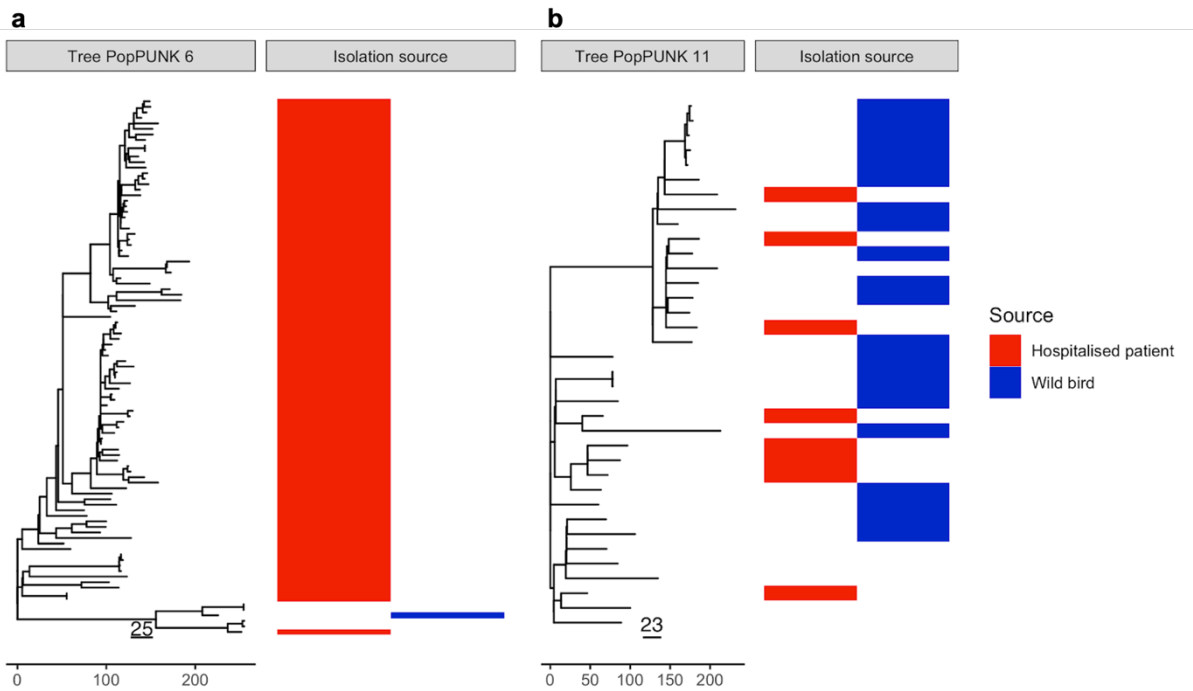
165

166

167

168

169



170

171

Supplementary Fig. 10: Hospital- and wild bird-associated clusters depict links between

172

the isolation sources within phylogenetic construction. An example of a hospital-associated

173

(HA) Population Partitioning Using Nucleotide K-mers (PopPUNK; PP)⁷ cluster PP6 ($n = 97$)

174

(**a**) and cluster PP11 ($n = 36$) (**b**), including isolates of both clinical (red) and wild bird (dark

175

blue) origin, aligned with a maximum-likelihood (ML) cluster phylogeny. Source data are

176

provided as a Source Data file.

177

178

179

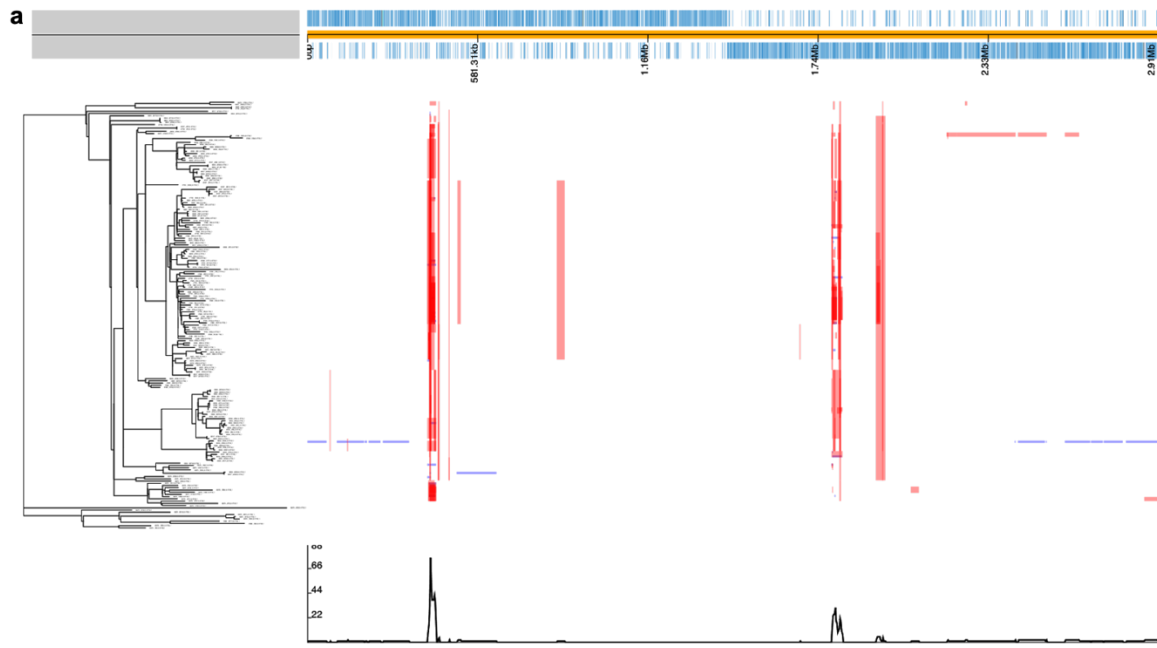
180

181

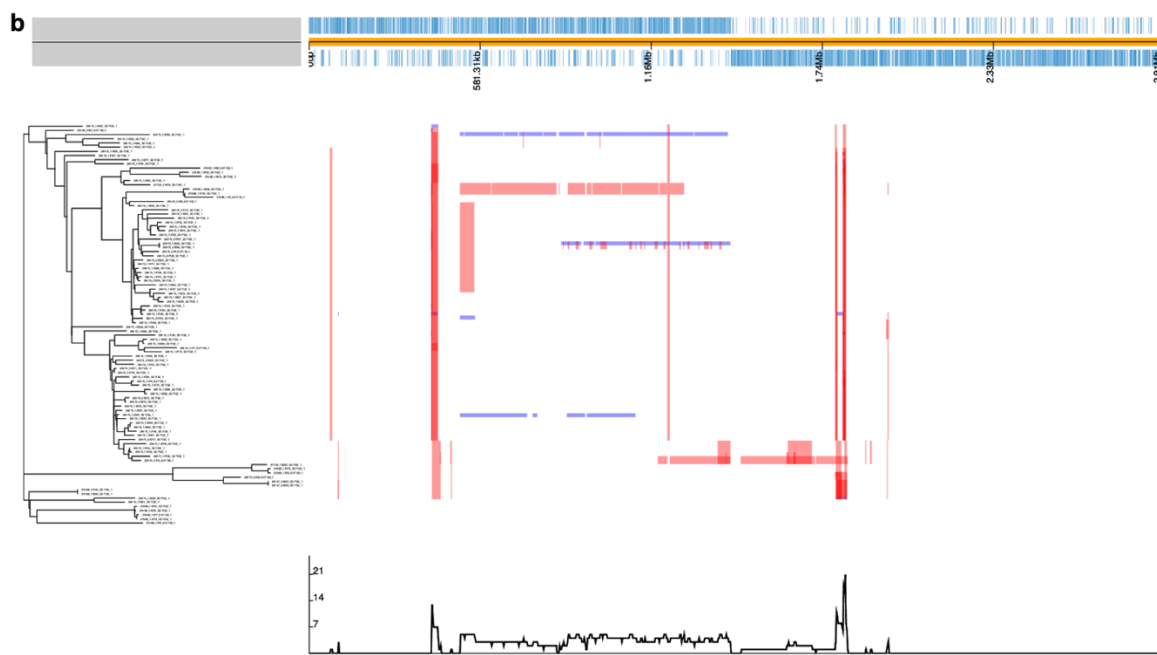
182

183

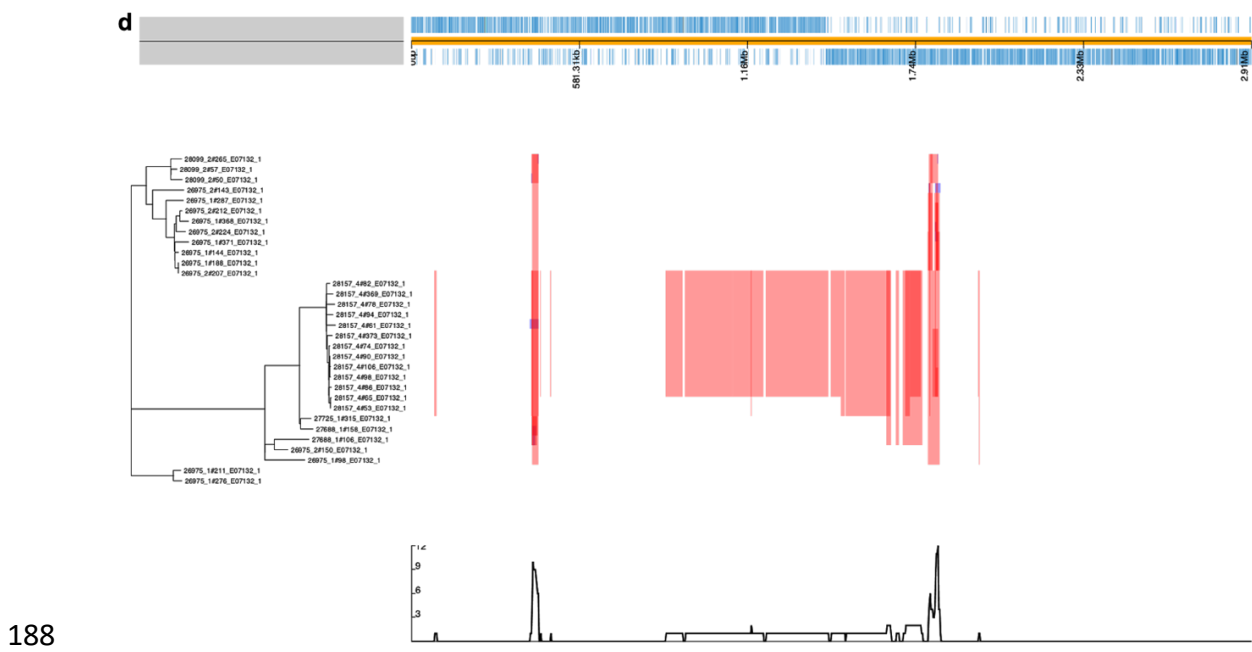
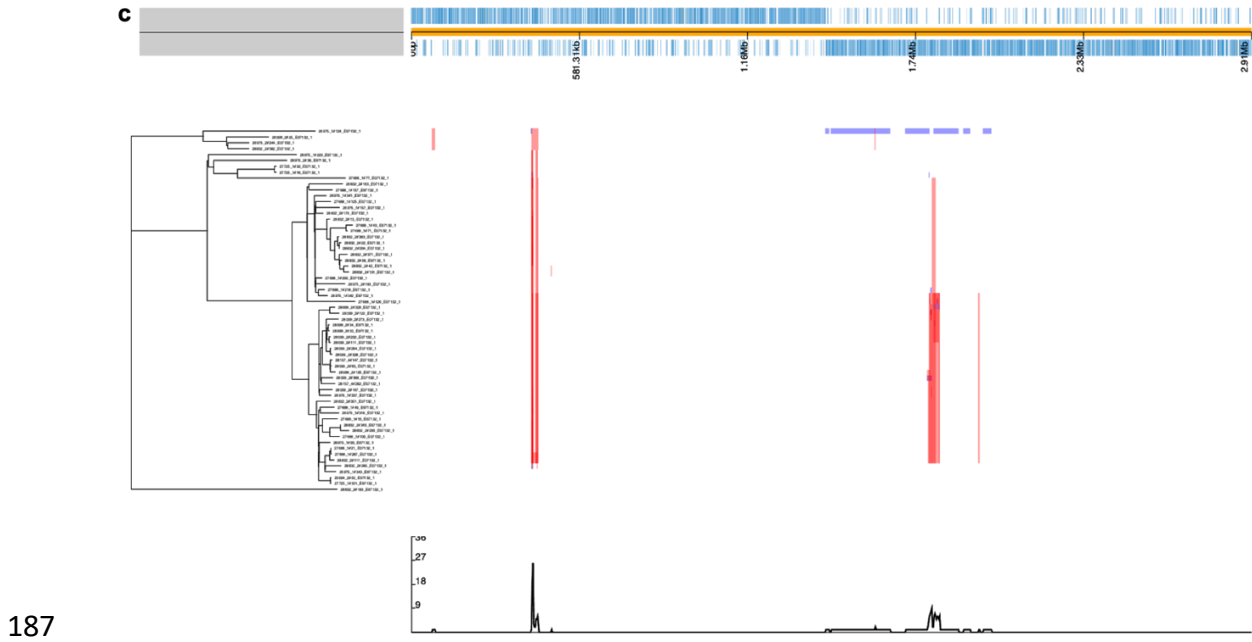
184



185



186



189 **Supplementary Fig. 11: Recombination predictions by using Genealogies Unbiased By**
 190 **recombinations In Nucleotide Sequences (Gubbins)⁹, for hospital-associated (HA)**
 191 **Population Partitioning Using Nucleotide K-mers (PopPUNK)⁷ clusters (PP). a, PP2 (*n* =**
 192 **193). b, PP6 (*n* = 97). c, PP7 (*n* = 62). d, PP18 (*n* = 32). Alignments on the hybrid assembly of**
 193 ***E. faecalis* old isolate E07132. Recombination blocks across maximum-likelihood (ML) cluster**
 194 **phylogeny are depicted by using Phandango¹⁰: red colouring indicates ancestral blocks**

195 (occurring at a non-terminal node), while blue indicates single isolates only. Reference
196 genome annotation panel above the recombination blocks shows the linearised genome,
197 with genes depicted as blue rectangles. Number of recombination events is plotted
198 underneath the recombination blocks.

199

200

201

202

203

204

205

206

207

208

209

210

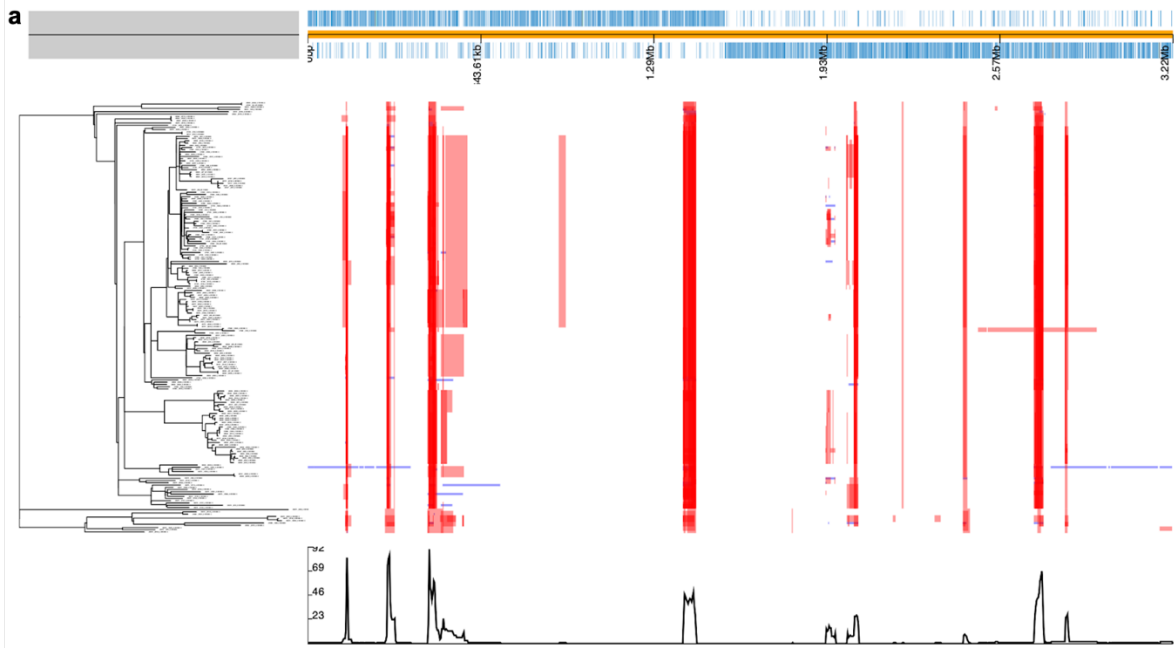
211

212

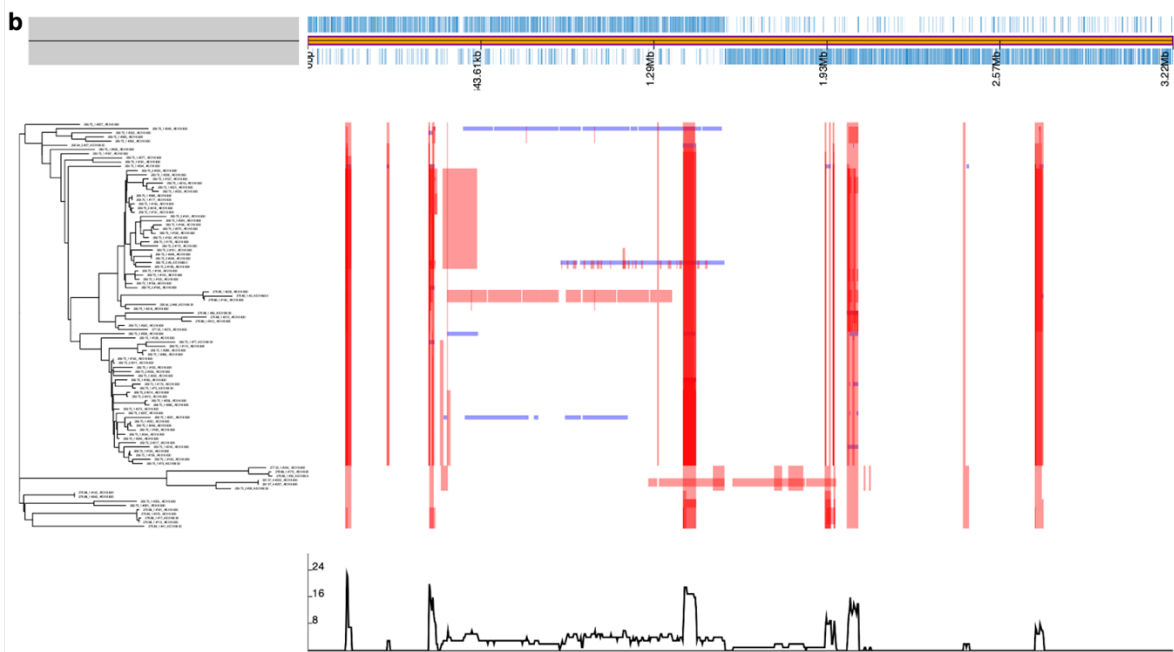
213

214

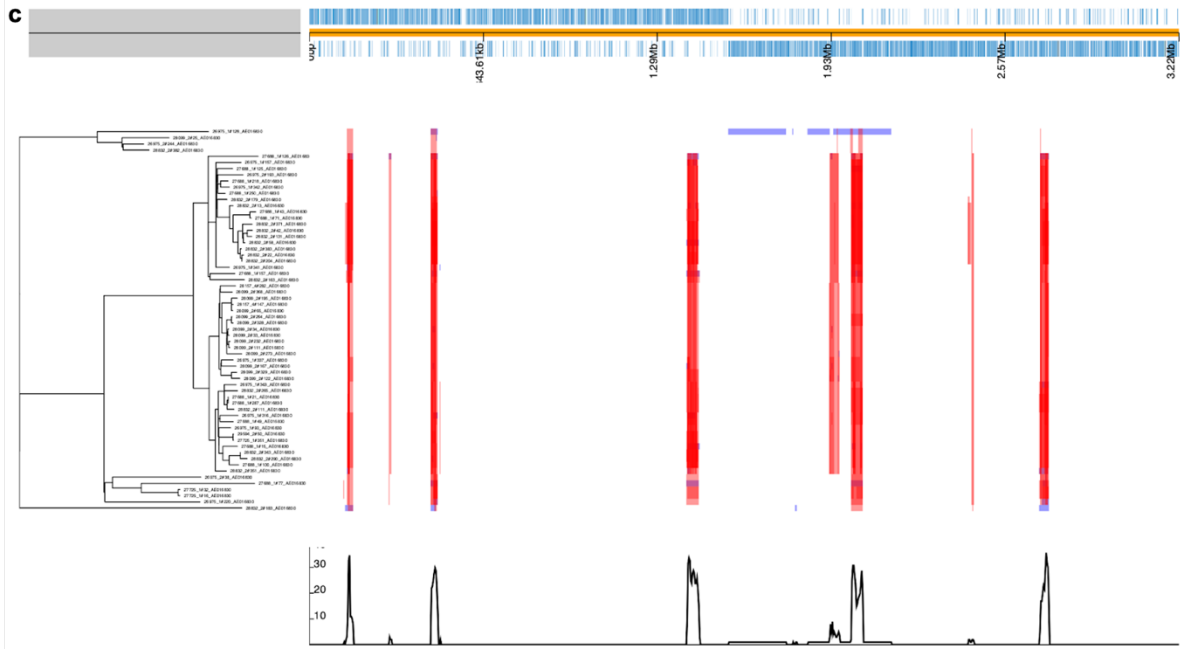
215



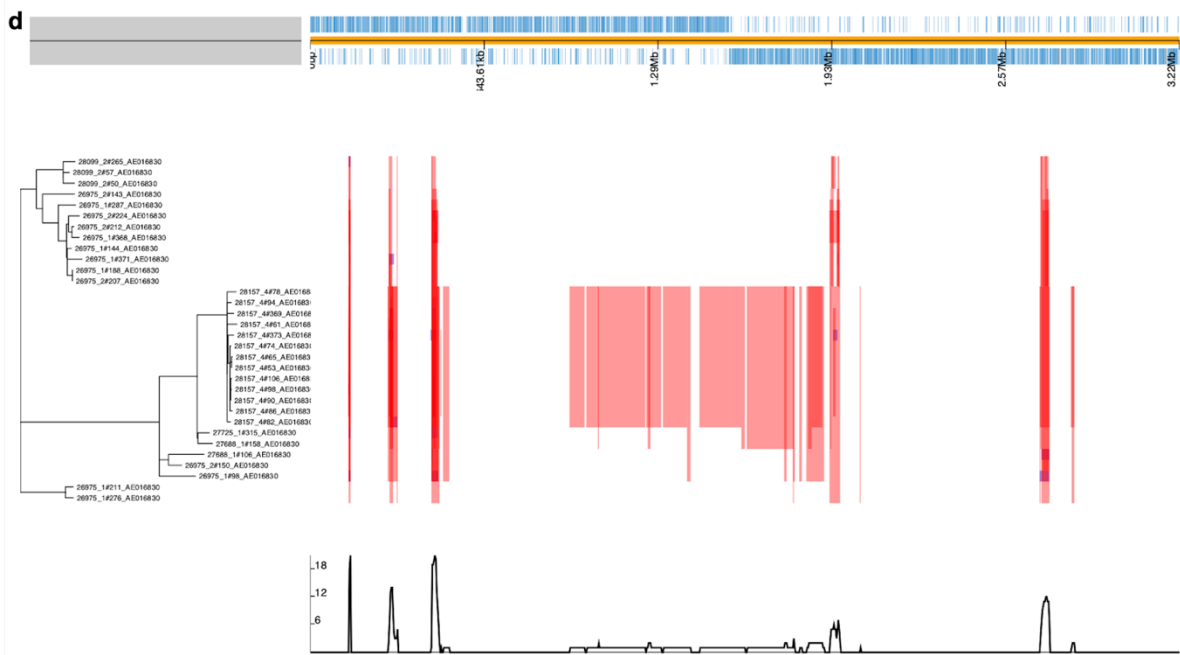
216



217



218



219 **Supplementary Fig. 12: Recombination predictions by using Genealogies Unbiased By**

220 **recomBinations In Nucleotide Sequences (Gubbins)⁹, for hospital-associated (HA)**

221 **Population Partitioning Using Nucleotide K-mers (PopPUNK)⁷ clusters (PP). a, PP2 (*n* =**

222 **193). b, PP6 (*n* = 97). c, PP7 (*n* = 62). d, PP18 (*n* = 32). Alignments on *E. faecalis* V583.**

223 Recombination blocks across maximum-likelihood (ML) cluster phylogeny are depicted by

224 using Phandango¹⁰: red colouring indicates ancestral blocks (occurring at a non-terminal

225 node), while blue indicates single isolates only. Reference genome annotation panel above
226 the recombination blocks shows the linearised genome, with genes depicted as blue
227 rectangles. Number of recombination events is plotted underneath the recombination
228 blocks.

229

230

231

232

233

234

235

236

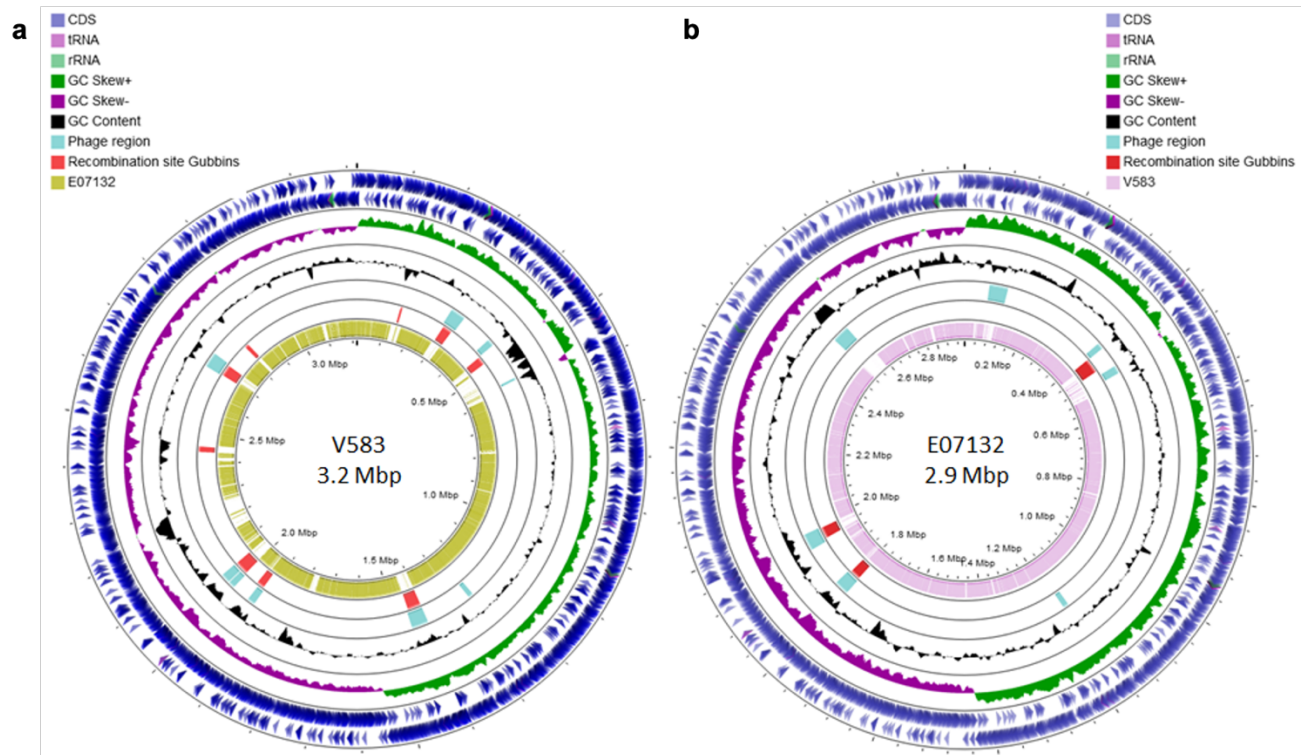
237

238

239

240

241



242

243 **Supplementary Fig. 13: Genomic maps of V583 (a) and E07132 (b) used as references in**

244 **the Genealogies Unbiased By recomBinations In Nucleotide Sequences (Gubbins)⁹**

245 **analysis.** Coloured as indicated in the legends, rings from outside to inside: CDS (blue), tRNA

246 (pink), and rRNA (light green) genes; GC Skew + (green)/- (purple); GC content (black); phage

247 integration sites as determined by Phaster¹¹ (light blue); recombination sites as determined

248 by Gubbins⁹ (red); and BLAST of the E07132 (ST97) against V583 (a) and V583 (ST6) against

249 E07132 (b). White blocks indicate absence of genes.

250

251

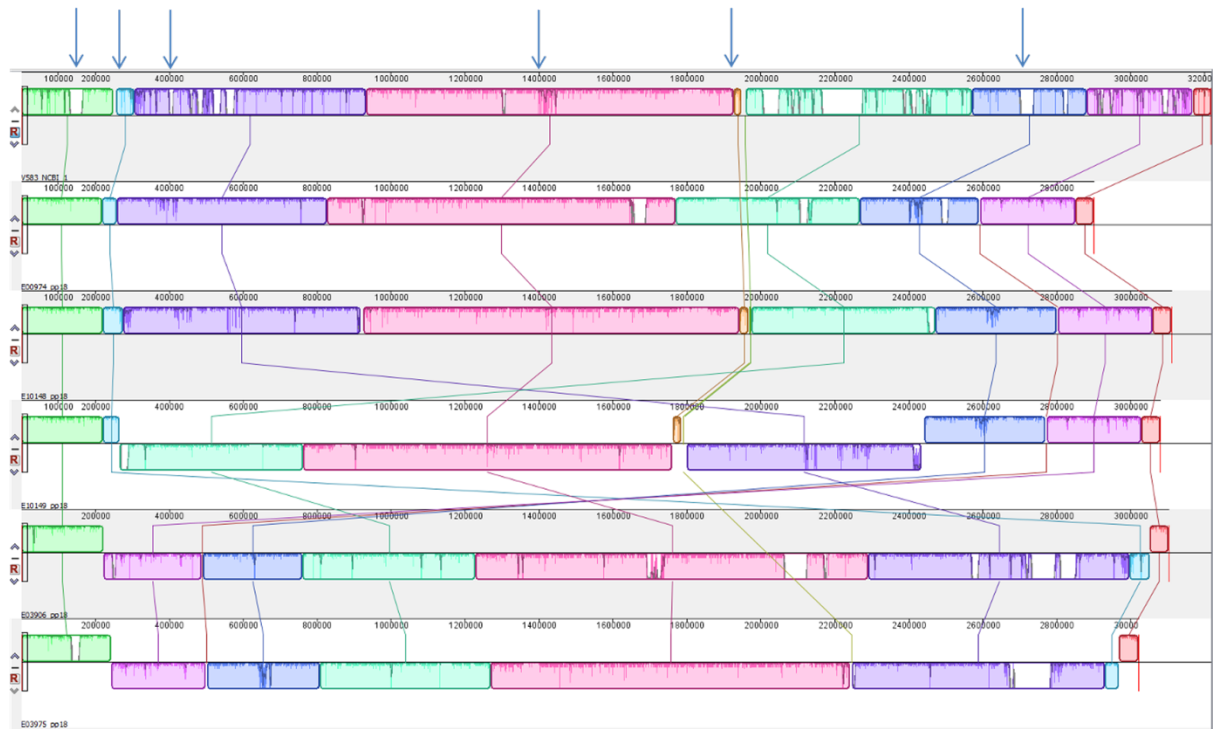
252

253

254

255

256



257

258 **Supplementary Fig. 14: MAUVE⁸ analysis on genome comparison of hospital-associated**

259 **(HA) PP18 isolates.** Coloured blocks indicating genomic rearrangements, and blue arrows

260 indicating recombination peaks as determined by Genealogies Unbiased By recomBinations

261 In Nucleotide Sequences (Gubbins)⁹.

262

263

264

265

266

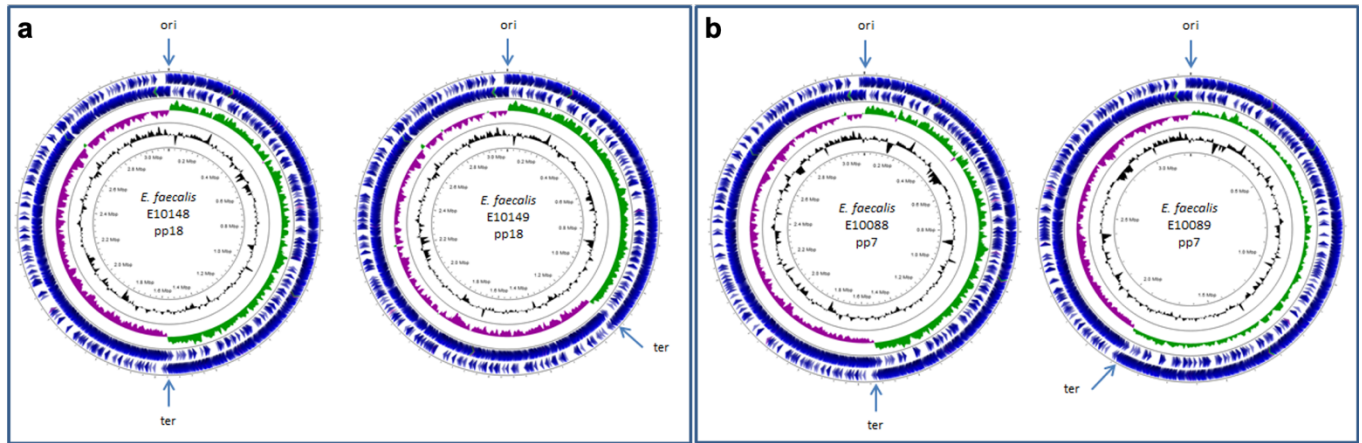
267

268

269

270

271



272

273 **Supplementary Fig. 15: Genomic maps of PP18 (panel a) and PP7 (panel b) strains**

274 **depicting a balanced replichore (left in panel) and replichore imbalance (right in panel).**

275 Rings from outside to inside: CDS (blue), tRNA (pink), and rRNA (light green) genes; GC Skew

276 + (green)/- (purple); GC Content (black). Arrows indicate the ori and ter sites.

277

278

279

280

281

282

283

284

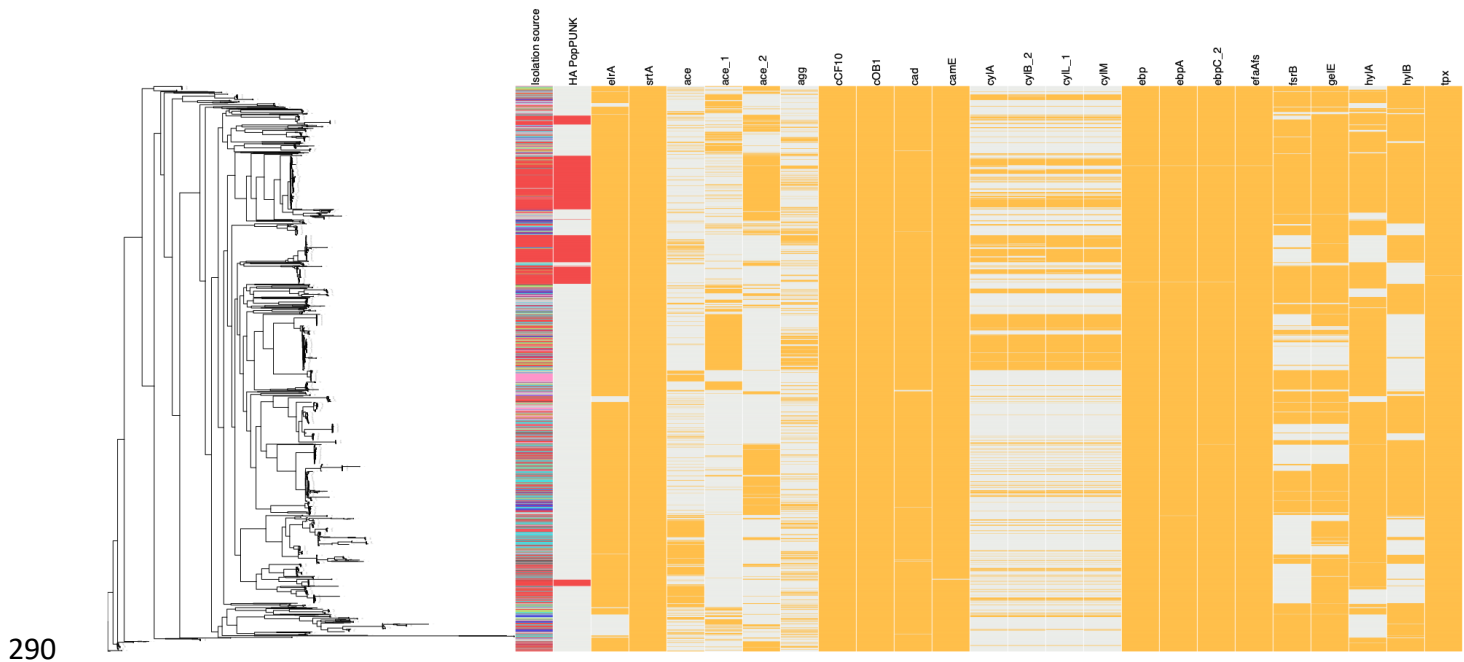
285

286

287

288

289



290

291 **Supplementary Fig. 16. Collection-wide virulence gene patterns of *E. faecalis* isolates ($n =$**
 292 **2,027) depict abundance of virulence genes but limited association with hospital isolates.**

293 First column, isolation sources: hospitalised patient (red), wild bird (dark blue), non-

294 hospitalised person (light blue), old human isolates (black), environment (green), farm

295 animal (pink), and others (grey). Second column, hospital-associated (HA) Population

296 Partitioning Using Nucleotide K-mers (PopPUNK)⁷ clusters (red). Presence (orange) and

297 absence (grey) of virulence genes, as defined by using Antimicrobial Resistance

298 Identification By Assembly (ARIBA)¹² against VirulenceFinder 2.0 database¹³, is aligned with

299 the species-wide reference mapping-based maximum likelihood (ML) phylogeny and

300 depicted by using Phandango¹⁰.

Supplementary Table 1. Analysing temporal signal and dating of hospital-associated (HA) Population Partitioning Using Nucleotide K-mers (PopPUNK) v.1.2.2⁷ clusters by using TempEst v.1.5.3¹⁴, Least-squares dating (LSD) v.0.3beta¹⁵, and Bayesian Evolutionary Analysis by Sampling Trees (BEAST2) v.2.5.0^{16–18}.

HA PopPUNK cluster	TempEst				LSD		BEAST2 ^a					
	tMRCA ^b	Slope (rate) ^c	Correlation coefficient	R ²	tMRCA ^b	Rate ^c	Constant tree			Exponential tree		
							tMRCA ^b	95% HPD ^d interval	clockRate ^e	tMRCA ^b	95% HPD ^d interval	clockRate ^e
PP2	1868	2.9391	0.5432	0.2950	1804	1.147	1846	1823–1866	5.719E-4	1844	1822–1865	5.635E-4
PP6	1987	9.4223	0.8512	0.7245	1965	4.200	1967	1961–1973	1.975E-3	1967	1961–1973	1.967E-3
PP7 ^f	1928	5.2709	0.5887	0.3466	-	-	-	-	-	-	-	-
PP18	1959	5.1903	0.9385	0.8808	1917	2.645	1917	1891–1941	2.847E-3	1921	1896–1943	2.948E-3
PP20 ^g	1637	1.1303	0.4487	0.2014	1542	0.958	NA	NA	NA	NA	NA	NA

^a Three replicate BEAST2 runs combined by using LogCombiner v.2.5.1.

^b tMRCA, the most recent common ancestor.

^c Rate (TempEst, LSD), estimate of the rate of evolution in substitutions per site per year.

^d HPD, highest posterior density.

^e clockRate (BEAST2), rate of evolution averaged over the whole tree and all sites.

^f PP7, TempEst on a subcluster of 55 isolates; LSD and BEAST2 on PP7 failed.

^g PP20 excluded from BEAST2 dating analyses due to poor temporal signal.

301

302

303

304

305

306

307 **References**

- 308 1. Argimón, S. *et al.* Microreact: visualizing and sharing data for genomic epidemiology
309 and phylogeography. *Microb Genom* **2**, e000093 (2016).
- 310 2. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo
311 pipeline. *Genome Biol.* **21**, 180 (2020).
- 312 3. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of
313 biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
- 314 4. Arredondo-Alonso, S. *et al.* Plasmids Shaped the Recent Emergence of the Major
315 Nosocomial Pathogen *Enterococcus faecium*. *MBio* **11**, (2020).
- 316 5. Arredondo-Alonso, S. *et al.* mlplasmids: a user-friendly tool to predict plasmid-and
317 chromosome-derived sequences for single species. *Microbial genomics* **4**, (2018).
- 318 6. Katz, L. S. *et al.* Mashtree: a rapid comparison of whole genome sequence files.
319 *Journal of Open Source Software* **4**, 1762 (2019).
- 320 7. Lees, J. A. *et al.* Fast and flexible bacterial genomic epidemiology with PopPUNK.
321 *Genome Res.* **29**, 304–316 (2019).
- 322 8. Darling, A. E., Mau, B. & Perna, N. T. progressiveMauve: multiple genome alignment
323 with gene gain, loss and rearrangement. *PLoS One* **5**, e11147 (2010).
- 324 9. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant
325 bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
- 326 10. Hadfield, J. *et al.* Phandango: an interactive viewer for bacterial population
327 genomics. *Bioinformatics* **34**, 292–293 (2018).
- 328 11. Arndt, D. *et al.* PHASTER: a better, faster version of the PHAST phage search tool.
329 *Nucleic Acids Res.* **44**, W16–21 (2016).

- 330 12. Hunt, M. *et al.* ARIBA: rapid antimicrobial resistance genotyping directly from
331 sequencing reads. *Microbial Genomics* vol. 3 (2017).
- 332 13. Joensen, K. G. *et al.* Real-time whole-genome sequencing for routine typing,
333 surveillance, and outbreak detection of verotoxigenic *Escherichia coli*. *J. Clin. Microbiol.* **52**,
334 1501–1510 (2014).
- 335 14. Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal
336 structure of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* **2**,
337 vew007 (2016).
- 338 15. To, T.-H., Jung, M., Lycett, S. & Gascuel, O. Fast Dating Using Least-Squares Criteria
339 and Algorithms. *Syst. Biol.* **65**, 82–97 (2016).
- 340 16. Drummond, A. J. & Rambaut, A. BEAST: Bayesian evolutionary analysis by sampling
341 trees. *BMC Evolutionary Biology* vol. 7 214 (2007).
- 342 17. Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis.
343 *PLoS Comput. Biol.* **10**, e1003537 (2014).
- 344 18. Bouckaert, R. *et al.* BEAST 2.5: An advanced software platform for Bayesian
345 evolutionary analysis. *PLoS Comput. Biol.* **15**, e1006650 (2019).