

1 **Genome-based characterization of hospital-adapted *Enterococcus faecalis* lineages**

2 Kathy E. Raven,^{1*} Sandra Reuter,¹ Theodore Gouliouris,^{1,2,3} Rosy Reynolds^{4,5}, Julie E Russell⁶,
3 Nicholas M. Brown^{2,4}, M. Estée Török,^{1,2,3} Julian Parkhill,⁷ Sharon J. Peacock.^{1,3,7,8}

4

5 ¹University of Cambridge, Department of Medicine, Box 157 Addenbrooke's Hospital, Hills
6 Road, Cambridge CB2 0QQ, United Kingdom

7 ²Public Health England, Clinical Microbiology and Public Health Laboratory, Box 236,
8 Addenbrooke's Hospital, Hills Road, Cambridge CB2 0QQ, United Kingdom

9 ³Cambridge University Hospitals NHS Foundation Trust, Hills Road, Cambridge CB2 0QQ,
10 United Kingdom

11 ⁴British Society for Antimicrobial Chemotherapy, Griffin House, 53 Regent Place Birmingham
12 B1 3NJ, United Kingdom

13 ⁵North Bristol NHS Trust, Southmead Hospital, Bristol, BS10 5NB, United Kingdom

14 ⁶Culture Collections, Public Health England, Porton Down, Salisbury SP4 0JG, United Kingdom

15 ⁷Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge,
16 United Kingdom

17 ⁸London School of Hygiene and Tropical Medicine, London, WC1E 7HT, United Kingdom

18

19

20 *Corresponding author: ker37@medschl.cam.ac.uk

21 **Vancomycin-resistant *Enterococcus faecalis* (VREfs) is an important hospital-**
22 **adapted pathogen. We undertook whole genome sequencing of *E. faecalis* associated with**
23 **bloodstream infection in the United Kingdom and Ireland (UK&I) over more than a decade**
24 **to determine the population structure and genetic associations with hospital adaptation.**
25 **Three lineages predominated in the population, two of which (L1 &L2) were nationally**
26 **distributed and L3 was geographically restricted. Genome comparison with a global**
27 **collection identified that L1 and L3 were also present in the United States but were**
28 **genetically distinct. Over 90% of VREfs belonged to L1-3, with resistance acquired and lost**
29 **multiple times in L1 and L2, but only once followed by clonal expansion in L3. Putative**
30 **virulence and antibiotic resistance genes were over-represented in L1, L2 and L3 isolates**
31 **combined, versus the remainder. Each of the 3 main lineages contained a mixture of**
32 **vancomycin-resistant and –susceptible *E. faecalis* (VSEfs), which has important implications**
33 **for infection control and antibiotic stewardship.**

34

35 Enterococci are the second and third most frequent cause of nosocomial infections in
36 the United States (US) and Europe respectively,^{1,2} with *Enterococcus faecalis* the most
37 commonly isolated species.² Vancomycin is the first-line antimicrobial drug for enterococci
38 with high-level resistance to ampicillin or for patients with penicillin allergy. Vancomycin
39 resistance was first reported in 1988³ and subsequently increased in prevalence. This rise was
40 predominantly due to *E. faecium*, but *E. faecalis* accounted for 11% of VRE bacteremias in the
41 UK&I between 2001-2013 (<http://www.bsacsurv.org>). Based on multi-locus sequence typing
42 (MLST)⁴, vancomycin resistance in *E. faecalis* has arisen in multiple genetic backgrounds,^{5,6}
43 and is associated with epidemic lineages⁷ Microbial genome sequencing provides the
44 opportunity to gain a detailed understanding of the molecular basis for hospital adaptation.
45 The first whole-genome sequence of *E. faecalis* was published in 2003.⁸ Subsequent genome
46 studies have compared 18 *E. faecalis* strains, which demonstrated the contribution of mobile

47 genetic elements to the diversity of the species;⁹ and 25 clinical and 7 non-clinical isolates,
48 which revealed that both were comparable in gene content.¹⁰ The ability to sequence large
49 bacterial collections means that the molecular epidemiology and gene content of epidemic
50 and sporadic lineages can now be systematically defined.

51 We sequenced 168 *E. faecalis* isolates (58 VREfs, 110 VSEfs) from national (British
52 Society for Antimicrobial Chemotherapy (BSAC) n=94), local (Cambridge University Hospitals
53 NHS Foundation Trust (Addenbrooke's Hospital and The Rosie Hospital) (CUH, n=60), and
54 reference (National Collection of Type Cultures (NCTC) n=14) collections (see Supplementary
55 Table 1). BSAC isolates originated from 21 UK&I hospitals between 2001 and 2011 (see
56 Supplementary Fig. 1 for geographical and temporal distribution), and CUH isolates were
57 collected between 2006 and 2012. All BSAC and CUH isolates were associated with
58 bloodstream infection. NCTC isolates were from humans, livestock and food products and
59 were predominantly isolated prior to 1951 (10/14 isolates).

60 Comparison of these genomes against the core genome of *E. faecalis* V583 identified
61 124,194 single nucleotide polymorphisms (SNPs) over 2,886,189 nucleotides (Fig. 1). A
62 striking feature of the phylogenetic tree based on these SNPs was that 53% of isolates
63 clustered into three distinct lineages (termed L1, L2 & L3). L1 was represented in each year of
64 the collection, while L2 was most frequently represented between 2001 and 2006 after which
65 there were only two isolates identified that belonged to this lineage (Supplementary Fig. 1).
66 This suggests clonal replacement, a phenomenon observed for other hospital-related
67 pathogens such as methicillin-resistant *Staphylococcus aureus*.¹¹ Annotation of lineage-
68 specific trees with geographical location demonstrated that L1 and L2 were nationally
69 distributed (epidemic) clones, whilst L3 was only isolated in two locations (CUH and a hospital
70 in the East Midlands referral network¹²) (Fig. 1). L3 was the dominant lineage at CUH (17/60
71 study isolates) and accounted for 3/13 isolates from a hospital in the East Midlands, the
72 phylogenetic tree supporting a single introduction into this hospital followed by local

73 diversification. We also explored the phylogenetic origins of VREfs by including all available *E.*
74 *faecalis* isolates from the NCTC collection. Eleven of 14 NCTC isolates clustered with recent
75 clinical isolates (Fig. 1), including seven isolated prior to 1951, and two VREfs from 1986 (the
76 first year that VREfs was recognised) that belonged to L1 or L2, which in the case of L2 may
77 represent a founder of the circulating vancomycin-resistant *E. faecalis* lineage.

78 To place our collection into a global context we compared these to the *E. faecalis*
79 genomes of isolates from around the world. This was achieved by retrieving all of the *E.*
80 *faecalis* genomes (n=353) held by the European Nucleotide Archive (ENA) as of 10/09/2015,
81 and combining our data with 347 of these (excluding 6 based on data quality). The
82 phylogenetic tree based on 1,293 genes conserved in 99% of these 515 isolates revealed that
83 isolates contained within L1 and L3 that originated from the UK and United States were
84 genetically distinct (Fig. 2). This indicates independent clonal expansion of dominant lineages
85 with limited international dissemination. One explanation for this is that these lineages are
86 hospital-associated with limited carriage beyond hospitals. Studies investigating community
87 carriage of VRE in the US and UK have failed to identify VREfs.^{13,14} By contrast, global isolates
88 from the non-dominant STs were closely related to UK isolates.

89 To compare our findings with those of published studies based on MLST, we assigned
90 STs to all 168 study isolates (see Supplementary Table 1). Isolates in L1 were assigned to ST6,
91 ST384 and ST642 (CC2), and L2 isolates were ST28 and ST640 (CC87). Both CCs have been
92 described as high-risk lineages based on their association with hospital-derived isolates in
93 Europe.⁵⁻⁷ L3 isolates were ST103 (CC388), which has only been reported previously in
94 relation to five clinical and two fecal isolates, all from the Americas.¹⁰

95 Comparison of the number of core genome SNPs for L1, L2 and L3 revealed lower
96 genetic diversity for L3 (range 3-60 SNPs, median 33 SNPs), compared with L1 (range 2-375
97 SNPs median 30 SNPs), and L2 (range 0-237 SNPs, median 139 SNPs). This led us to use
98 Bayesian Evolutionary Analysis Sampling Trees (BEAST) to date these lineages.¹⁵ The last

99 common ancestor of L3 was estimated to be 1998 (95% highest posterior density (HPD)
100 interval, 1980-2004) (Supplementary Fig. 2a), consistent with the earliest reported isolation
101 of ST103 in the literature of 2002.¹⁰ The last common ancestor of L1 was predicted to be 1918
102 (95% HPD interval, 1868-1960), with a clonal expansion in 1997 (95% HPD interval, 1992-
103 2000) (Supplementary Fig. 2b). The early estimate for the last common ancestor of L1 relied
104 on just three outlying isolates, and so a second algorithm was used to detect and remove
105 recombination events and the BEAST analysis repeated to rule out the role of undetected
106 recombination. This predicted a last common ancestor in 1852 (95% HPD interval, 1811-
107 1956). It has been proposed previously that CC2 (L1) emerged recently, based on the lack of
108 isolates identified prior to the 1980s.¹⁶ Our analysis indicates that this lineage may have been
109 in existence since the mid 1850s to early 1900s, with a clonal expansion in the 1990s. BEAST
110 analysis of L2 failed, probably because of a limited number of isolates with high genetic
111 diversity and wide temporal spread.

112 Establishing the rate of mutation in the core genome provides a molecular clock that
113 contextualises analyses of bacterial genomes during putative outbreak investigations,¹⁷⁻¹⁹ but
114 has not been defined previously for *E. faecalis*. The rate of evolution was estimated to be
115 8.18×10^{-7} SNPs/site/year (approximately 2.5 SNPs/year) for L1 and 1.14×10^{-6} SNPs/site/year
116 (approximately 3.4 SNPs/year) for L3. Based on these mutation rates and patient ward
117 locations, we excluded direct patient-to-patient transmission of the CUH study isolates.

118 We explored the genetic basis for the success of the dominant *E. faecalis* lineages
119 using a candidate gene approach by comparing the prevalence of putative virulence and
120 antibiotic resistance genes in L1, L2 and L3 isolates combined, versus the remainder. *ace*,
121 *gelE*, *asa1*, *agg*, *cyl*, *elrA*, and genes conferring resistance to tetracyclines, aminoglycosides,
122 trimethoprim, chloramphenicol, macrolides/lincosamides/streptogramin B (MLSB),
123 quaternary ammonium compounds (qacs) and vancomycin were over-represented in L1-3
124 compared to the rest (Fig. 3). There was a striking difference in the prevalence of genes

125 encoding aminoglycoside and vancomycin resistance, two commonly used antibiotics for
126 enterococcal infection, in dominant versus non-dominant lineages. Our findings extend
127 previous reports that epidemic lineages are enriched for multi-drug resistance and specific
128 virulence determinants.^{6,7} We then compared the prevalence of the candidate virulence
129 genes in VREfs versus VSEfs contained in L1, L2 and L3 (Supplementary Fig. 3). This showed no
130 significant difference, indicating that over-representation of virulence genes is lineage- rather
131 than VRE-specific.

132 We then analysed the pangenome²⁰ of the 168 isolates to obtain a more detailed
133 understanding of their entire genomic repertoire. This indicated that *E. faecalis* has an open
134 genome with a gamma value of 0.21 (Supplementary Fig. 4), corroborating results derived
135 previously from the analysis of 5 genomes.²¹ The pangenome contained 8,202 genes, of
136 which 1,967 genes were conserved across the collection. Of the 6,235 genes in the accessory
137 genome, 1,687 were present just once. The most common accessory genes encoded
138 hypothetical proteins (n=2,558), IS elements or transposons (n=177), phage or plasmid-
139 associated proteins (n=462 and n=113 respectively), transcriptional regulators (n=225), ABC
140 transporters or cassettes (n=124), and phosphotransferase systems (n=118). A total of 819
141 genes were only found in the three dominant lineages, of which 109 were present in more
142 than 10 isolates (Supplementary Table 2), including a WxL domain surface protein unique to
143 L1. Comparison of the amino acid sequence of the WxL protein from L1 to the proteome of
144 V583 revealed a 100% match to EF_3248, one of 27 WxL proteins identified by Brinster *et*
145 *al.*²² No genes or homoplasic non-synonymous SNPs were ubiquitous in the dominant
146 lineages and absent from all sporadic lineages, suggesting that there is no single factor that
147 contributed to the emergence of these dominant clones, although antibiotic resistance and
148 virulence determinants are likely to represent multifactorial contributory factors. Analysis of
149 non-synonymous SNPs unique to L3 revealed 122 SNPs in 95 genes (Supplementary Table 3),

150 but no single genetic event was identified that might explain the geographically constrained
151 success of this lineage.

152 Recombination is thought to be a major mechanism by which the *E. faecalis* genome
153 evolves, which led us to estimate sites of recombination in the core genome²³ for L1, L2 and
154 L3 (Supplementary Fig. 5). Recombination accounted for 12.3% of the core genome in L1
155 (6.5% related to a large recombination event in 2 isolates) and 3.9% in L2, with a single
156 predicted 4 bp recombination event in one L3 isolate. This contrasts with reports that
157 recombination across the species is high,⁴ which led us to use an alternative algorithm
158 (BratNextGen²⁴) to detect recombination. This revealed similarly low levels of recombination
159 in L2 and L3 (6.2% and 0.3% respectively) but higher rates in L1 (37%), although most of this
160 (93%) was contained within two large recombination events (Supplementary Fig. 5). One
161 possible explanation for the low levels of recombination is that this drove the initial
162 diversification of the species but subsequently contributed little to short-term evolution.

163 Finally, we analysed the genetic basis of vancomycin resistance in the collection.
164 Nearly all VREfs (57/58) carried *vanA*, with a single NCTC isolate carrying *vanB*. Annotation of
165 the tree with resistance to vancomycin showed that all three dominant lineages contained a
166 mixture of VREfs and VSEfs, with 89% of BSAC and 95% of CUH VREfs belonging to L1-3. Based
167 on mapping to a reference *vanA* transposon (Tn1546) there was no SNP-based variation
168 between transposons with the exception of one that had a C → T substitution at position
169 5745. However, there was substantial variation in gene content. The transposase, resolvase,
170 *vanY* and *vanZ* genes were not detected in some isolates, but despite this the minimum
171 inhibitory concentration (available for the 35 *vanA* positive BSAC isolates) was consistently
172 very high (≥256 mg/L). There was considerable variation in genetic content within and
173 between the L1 and L2 transposon, whilst L3 had limited variation with two variants relating
174 to VREfs isolated in 2006-2009 and 2009-2012 respectively, and three partial deletions (Fig.
175 4). Analysis of the insertion sites for Tn1546 revealed multiple insertion sites for L1 and L2,

176 but only one site was identified for L3 in the 11/14 genomes for which this analysis proved
177 possible (Fig. 4 and Supplementary Table 4). Analysis using BLAST revealed that these
178 insertion sites were best matched to plasmids, a finding corroborated using plasmid
179 extraction and *vanA* hybridization for insertion site types 1A, 1B, 2B and 3 (data not shown).
180 These data indicate multiple acquisition and loss of the *vanA* transposon in L1 and L2,
181 suggesting a significant fitness cost. By contrast, the single acquisition followed by clonal
182 expansion in L3 suggests that the transposon has negligible cost or confers a benefit in this
183 lineage. Foucault *et al.*²⁵ demonstrated that its integration site in the chromosome
184 predominantly determined the fitness cost of the *vanB* transposon. One possible reason for
185 the retention of *vanA* in L3 is that the transposon has inserted into the plasmid at a location
186 that lacks a fitness cost to the bacterium. However, *vanA* is inserted at the same site in 10
187 isolates from L1 and L2 and there is limited evidence for retention of *vanA* in these isolates.

188 In conclusion, whole genome sequencing of *E. faecalis* has highlighted the dominance
189 of epidemic lineages in the UK&I, but also showed that a lineage with features of an epidemic
190 lineage was confined to two hospitals. Additionally, we identified that the UK and US have
191 genetically distinct populations belonging to two of these lineages, suggesting a lack of
192 international transmission. The mutation rate defined here will have utility in clinical practice
193 as sequencing technology is introduced into the investigation of putative outbreaks. Genome-
194 level data provided comprehensive insights into the gene content of dominant versus
195 sporadic lineages and allowed us to describe the evolution of vancomycin resistance in this
196 collection, which included multiple loss and acquisition events. The observation that the
197 major VREfs lineages were also the common lineages for VSEfs has important implications for
198 infection control and antibiotic stewardship, since the control of VREfs is likely to depend on
199 defining and addressing drivers for VSEfs and its transmission.

200

201

202 **ONLINE METHODS**

203 **Ethical approval**

204 The study was approved by the National Research Ethics Service (ref: 12/EE/0439) and the
205 Cambridge University Hospitals NHS Foundation Trust (CUH) Research and Development
206 (R&D) Department.

207

208 **Isolate collection**

209 The 168 *Enterococcus faecalis* isolates used in this study were selected from three
210 collections: NCTC (n=14, deposited between 1927 and 2007), BSAC (n=94, isolated between
211 2001 and 2011) and CUH (n=60, isolated between Nov 2006 and Dec 2012). The collection
212 was enriched for vancomycin-resistant isolates by selecting all of the available VREfs from
213 NCTC (n=3) and BSAC (n=35) and the first stored isolate from all cases of VREfs bacteremia at
214 CUH (n=20). To relate this to the underlying VSEfs population, 110 VSEfs were selected as
215 follows: (i) all available VSEfs from the NCTC (n=11), (ii) 59 VSEfs from BSAC (35 matched to
216 the BSAC VREfs cases by hospital and year of isolation where available, and an additional 24
217 VSEfs to gain greater representation of the VSEfs population); (iii) 40 VSEfs from CUH (the
218 first stored bacteremia-associated isolate matched to CUH VREfs cases by isolation date
219 (n=17), or that occurred 30 days or more after admission (n=19), and 4 additional VSEfs that
220 were available to increase the representation of the local population). BSAC hospitals were
221 assigned to referral networks described previously,¹² which are clusters of hospitals more
222 likely to exchange patients within the cluster than outside of that cluster.

223

224 **Microbiology and sequencing**

225 Bacterial isolates were cultured on Columbia Blood Agar (Oxoid, Basingstoke, UK) and
226 incubated at 37°C for 48 hours in air. Vancomycin susceptibility was determined using the
227 agar dilution method²⁶ (BSAC isolates), or the Vitek2 instrument (Biomérieux, Marcy l'Etoile,

228 France) with the AST-P607 card (CUH and NCTC VREfs isolates). DNA was extracted using the
229 QIAextractor (QIAGEN), according to the manufacturer's instructions. Library preparation was
230 conducted according to the Illumina protocol, and sequencing was performed on an Illumina
231 HiSeq2000 with 100-cycle paired-end runs. Sequence data for all isolates have been
232 submitted to the European Nucleotide Archive (ENA) (www.ebi.ac.uk/ena) with the accession
233 numbers shown in Supplementary Table 1.

234

235 **Phylogenetic analyses**

236 Sequence reads were mapped using SMALT
237 (<http://www.sanger.ac.uk/resources/software/smalt/>) to the *E. faecalis* reference genome
238 V583 (ENA accession number AE016830) for collection-wide analysis. This reference was
239 selected because it is one of only two finished *E. faecalis* genomes from clinical isolates, and
240 has been used in multiple studies (the second complete genome having only been published
241 in 2014). For analysis of lineages L1, L2 and L3, the oldest isolate from each lineage was
242 selected as a reference for mapping and an assembly created using Velvet. Mobile genetic
243 elements were identified using gene annotation, PHAST²⁷ (phast.wishartlab.com), WebACT²⁸
244 (<http://www.webact.org>) and BLAST²⁹ (blast.ncbi.nlm.nih.gov) and were excluded in addition
245 to contigs less than 500bp in length to create a 'core' genome. The core genome sizes were
246 2,886,189, 2,698,500bp, 2,372,434bp and 2,707,007bp for V583, L1, L2 and L3, respectively.
247 Single nucleotide polymorphisms (SNPs) in the core genome were determined using an in-
248 house script and used to estimate maximum likelihood trees using RAxML³⁰ with 100
249 bootstraps. Recombination was removed from the lineage-specific analyses using Gubbins.²³
250 To place the isolates into a global context, all of the available *E. faecalis* sequences listed in
251 GenBank were downloaded from the ENA (n=353). Six isolates were excluded due to poor
252 assemblies/annotation. The assemblies of the remaining 347 isolates were combined with
253 the assemblies of the study isolates (created using Velvet), annotated with Prokka, and a pan-

254 genome estimated using Roary.²⁰ A 90% identity cut-off was used and core genes were
255 defined as those in 99% of isolates. A maximum likelihood tree of the 25,294 SNPs in the
256 1,416 core genes was created using RAxML and 100 bootstraps. iTOL³¹ and FigTree were used
257 to visualise the trees. Assemblies were compared to the MLST database
258 (pubmlst.org/efaecalis/) sited at the University of Oxford³² using an in-house script.

259

260 **Population history and mutation rate**

261 Genetic diversity was calculated based on pairwise SNP differences. Bayesian Evolutionary
262 Analysis Sampling Trees (BEAST)¹⁵ was used to date the phylogeny and estimate a mutation
263 rate for L1 and L3 using the core genome after removal of regions of recombination using
264 Gubbins. BratNextGen²⁴ was used to verify the results of Gubbins for L1 using the following
265 parameters: 10 iterations, 100 permutation runs and a significance threshold of 0.05. One
266 NCTC isolate was excluded from the analysis for L1 because the isolation date was unknown.
267 A Hasegawa, Kishino and Yano (HKY) model and gamma distribution was used, and the best
268 molecular clock and tree selected based on Bayes factors calculated from path sampling and
269 stepping stone sampling^{33,34}: for L3 an exponential clock and constant tree were used, for L1 a
270 lognormal clock and Bayesian skyline tree were used, and for the repeat analysis of L1 (with
271 recombination events identified and removed based on BratNextGen) a lognormal clock and
272 constant tree were used.

273

274 **Detection of candidate genes**

275 Virulence genes were chosen based on evidence from experimental mammalian models³⁵ and
276 their presence determined by *in silico* PCR using previously published primers: *ace*,³⁶ *esp*,³⁷
277 *gelE*,³⁸ *asa1*,³⁹ *agg*³⁶ and *cyl*,⁴⁰ *elrA* (*OEF2* and *OEF8*)⁴¹, *gls24*,⁴² *tpx* (*ef1933for* and *tpxrev*)⁴³,
278 *bgsA* (*bgsA for* and *bgsArev*)⁴⁴, *srtA* (*EF3056F* and *EF3056R*)⁴⁵, *sigV* (*SVRT1-2*)⁴⁶, *epaA*
279 (*AB270_epa_F* and *AB271_epa_R*)⁴⁷, *epaB* (*AB272_epaB_F* and *AB273_epaB_R*)⁴⁷, *epaE*

280 (*AB276_epaE_F* and *AB277_epaE_R*)⁴⁷, *epaN* (*AB288_epaN_F* and *AN289_epaN_R*)⁴⁷ and
281 *perA* (*perA-FF* and *perA-RR*)⁴⁸. The presence of *msrA* and *msrB* were determined by coverage
282 of EF1681 and EF3164 respectively, when mapped to the V583 reference genome. The
283 presence of *vanA* and *vanB* were established by *in silico* PCR using published primers.^{35,49}
284 Genes encoding resistance to additional antimicrobial drugs were detected by comparing the
285 whole genome of each isolate with the ResFinder database (compiled in 2012),⁵⁰ which has
286 been manually curated since publication. Sequences were compared using an in-house script,
287 and genes with 100% match to length and > 90% identity match were classified as present. *In*
288 *silico* PCR using previously published primers was used for genes not in the ResFinder
289 database: *dfrF*⁵¹ and *qacZ*⁵². Statistical significance was determined using Fisher's exact test.

290

291 **Pangenome and recombination**

292 The pan genome was estimated using Roary.²⁰ Core genes were defined as those present in
293 all 168 isolates with a 90% ID cut-off. The proteome of *E. faecalis* strain V583 was
294 downloaded from the ENA and interrogated with the WxL protein described in this study
295 using the protein version of BLAST. Recombination was identified within L1-L3 using Gubbins
296 and verified using BratNextGen as described above.

297

298 **Characterising the Tn1546 transposon**

299 Sequence reads for each isolate were mapped to Tn1546 (accession number M97297) from
300 the *E. faecium* strain BM4147 using SMALT. The depth of coverage was between ~30x and
301 ~500x for all isolates, with 53/57 (93%) at a depth of > ~50x. To identify the insertion sites,
302 the sequences adjacent to the start and end of the Tn1546 transposons were extracted from
303 the assemblies up to a maximum of 10,000bp or the end of the contig. The sequences
304 adjacent to the start of Tn1546 were too short to analyse but sequences adjacent to the end
305 of Tn1546 (termed "insertion site sequences") were compared between isolates. Insertion

306 site sequences that were identical for more than 200bp were grouped (groups 1-3 in
307 Supplementary Table 4), and then subgroups defined if there were any changes in the
308 downstream sequence, with no evidence that an insertion or deletion explain this change
309 (changes described in Supplementary Table 4). Where the insertion site sequence available
310 was too short to determine which subgroup it belonged to, this was categorised into
311 subgroup A (the most prevalent subgroup) for simplification of Fig. 4. Each subgroup was
312 identified as plasmid or chromosome-based using BLAST, with transposons considered
313 plasmid-borne if the highest match was to a plasmid and there was no match above 25%
314 coverage to an *E. faecalis* chromosome. Insertion site sequences less than 250bp were not
315 considered long enough for an accurate identification. To verify whether the *vanA*
316 transposons were located on plasmids, plasmid extraction followed by *vanA* hybridization
317 was performed. Representative isolates were selected for each of the transposon insertion
318 sites defined using the sequence data. Plasmids were extracted using the Kado and Liu⁵³
319 method except that 100 mg/ml lysozyme was added with the E buffer followed by incubation
320 for 1 hour. Extracts were run on a 0.7% agarose gel and blotted using capillary transfer onto
321 Hybond N+ (Amersham, Buckinghamshire, UK). Hybridization was performed using the DIG-
322 high prime labeling and detection starter kit I (Roche Applied Science, Mannheim, Germany),
323 and luminescence was detected using CSPD (Roche Applied Science). *E. faecium* BM4147 and
324 NCTC 8132 were used as positive and negative controls, respectively, and two isolates with
325 known plasmid sizes were used as size markers (*Yersinia enterocolitica* YE212/92 (BT 2, O:9) and
326 YE53/03 (BT 1A, O:5)). Probes were created using the following primers: *vanA*-1: 5'-
327 GGGAAAACGACAATTGC-3', *vanA*-2: 5'GTACAATGCGGCCGTTA-3'.⁴⁹

328

329 **Accession codes**

330 The sequence data for the study isolates has been deposited in the ENA under the study
331 accessions PRJEB4344, PRJEB4345 and PRJEB4346, with the accession numbers for individual

332 isolates listed in Supplementary Table 1. Additional sequences used in this study were the *E.*
333 *faecalis* reference genome V583 (ENA accession number AE016830) and Tn1546 from the *E.*
334 *faecium* strain BM4147 (ENA accession number M97297).

335 **References**

- 336 1. Sievert, D. M. *et al.* Antimicrobial-resistant pathogens associated with healthcare-
337 associated infections: summary of data reported to the National Healthcare Safety
338 Network at the Centers for Disease Control and Prevention, 2009-2010. *Infect. Control*
339 *Hosp. Epidemiol.* **34**, 1–14 (2013).
- 340 2. European Centre for Disease Prevention and Control. Point prevalence survey of
341 healthcare-associated infections and antimicrobial use in the European acute care
342 hospitals. *ECDC* doi 10.2900/86011 (2013).
- 343 3. Uttley, A. H. C., Collins, C.H., Naidoo, J. & George, R. C. Vancomycin-resistant
344 enterococci. *Lancet* **2**, 57–58 (1988).
- 345 4. Ruiz-Garbajosa, P. *et al.* Multilocus sequence typing scheme for *Enterococcus faecalis*
346 reveals hospital-adapted genetic complexes in a background of high rates of
347 recombination. *J. Clin. Microbiol.* **44**, 2220–8 (2006).
- 348 5. Freitas, A. R., Novais, C., Ruiz-Garbajosa, P., Coque, T. M. & Peixe, L. Clonal expansion
349 within clonal complex 2 and spread of vancomycin-resistant plasmids among different
350 genetic lineages of *Enterococcus faecalis* from Portugal. *J. Antimicrob. Chemother.* **63**,
351 1104–11 (2009).
- 352 6. Kuch, A. *et al.* Insight into antimicrobial susceptibility and population structure of
353 contemporary human *Enterococcus faecalis* isolates from Europe. *J. Antimicrob.*
354 *Chemother.* **67**, 551–8 (2012).
- 355 7. Kawalec, M. *et al.* Clonal structure of *Enterococcus faecalis* isolated from Polish
356 hospitals: characterization of epidemic clones. *J. Clin. Microbiol.* **45**, 147–53 (2007).
- 357 8. Paulsen, I. T. *et al.* Role of mobile DNA in the evolution of vancomycin-resistant
358 *Enterococcus faecalis*. *Science* **299**, 2071–2074 (2003).
- 359 9. Palmer, K. L., *et al.* Comparative Genomics of Enterococci: Variation in *Enterococcus*
360 *faecalis*, Clade Structure in *E. faecium*, and Defining Characteristics of *E. gallinarum*
361 and *E. casseliflavus*. *MBio* **3**, 1–11 (2012).
- 362 10. Kim, E. B. & Marco, M. L. Nonclinical and Clinical *Enterococcus faecium* Strains, but
363 Not *Enterococcus faecalis* Strains, Have Distinct Structural and Functional Genomic
364 Features. *Appl. Environ. Microbiol.* **80**, 154–165 (2014).
- 365 11. Hsu, L.-Y. *et al.* Evolutionary dynamics of methicillin-resistant *Staphylococcus aureus*
366 within a healthcare system. *Genome Biol.* **16**, 81 (2015).
- 367 12. Donker, T., Wallinga, J., Slack, R. & Grundmann, H. Hospital networks and the dispersal
368 of hospital-acquired pathogens by patient transfer. *PLoS One* **7**, e35002 (2012).
- 369 13. Coque, T. M., Tomayko, J. F., Ricke, S. C., Okhyusen, P. C. & Murray, B. E. Vancomycin-
370 Resistant Enterococci from Nosocomial, Community, and Animal Sources in the United
371 States. *Antimicrob. Agents Chemother.* **40**, 2605–2609 (1996).

- 372 14. Jordens, J. Z., Bates, J. & Griffiths, D. T. Faecal carriage and nosocomial spread of
373 vancomycin-resistant *Enterococcus faecium*. *J. Antimicrob. Chemother.* **34**, 515–528
374 (1994).
- 375 15. Drummond, A. J., Suchard, M. A, Xie, D. & Rambaut, A. Bayesian phylogenetics with
376 BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–73 (2012).
- 377 16. Palmer, K. L. *et al.* Enterococcal Genomics. In Gilmore, M. S. *et al.*, editors.
378 *Enterococci: From Commensals to Leading Causes of Drug Resistant Infection*
379 [Internet]. Boston: Massachusetts Eye and Ear Infirmary (2014). Available from:
380 <http://www.ncbi.nlm.nih.gov/books/NBK190425/>
- 381 17. Köser, C. U. *et al.* Rapid Whole-Genome Sequencing for Investigation of a Neonatal
382 MRSA Outbreak. *N. Engl. J. Med.* **366**, 2267–2275 (2013).
- 383 18. Harris, S. R. *et al.* Evolution of MRSA During Hospital Transmission and
384 Intercontinental Spread. *Science* **327**, 469–474 (2010).
- 385 19. Walker, T. M. *et al.* Whole-genome sequencing to delineate *Mycobacterium*
386 *tuberculosis* outbreaks: a retrospective observational study. *Lancet Infect. Dis.* **13**,
387 137–146 (2013).
- 388 20. Page, A. J. *et al.* Roary: Rapid large-scale prokaryote pan genome analysis.
389 *Bioinformatics* pii: btv421 [Epub ahead of print] (2015).
- 390 21. The Human Microbiome Jumpstart Reference Strains Consortium. A Catalog of
391 Reference Genomes from the Human Microbiome. *Science* **328**, 994–999 (2010).
- 392 22. Brinster, S., Furlan, S. & Serror, P. C-Terminal WxL Domain Mediates Cell Wall Binding
393 in *Enterococcus faecalis* and Other Gram-Positive Bacteria. *J. Bacteriol.* **189**, 1244–
394 1253 (2007).
- 395 23. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant
396 bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2014).
- 397 24. Marttinen, P. *et al.* Detection of recombination events in bacterial genomes from large
398 population samples. *Nucleic Acids Res.* **40**, e6 (2012).
- 399 25. Foucault, M., Depardieu, F., Courvalin, P. & Grillot-Courvalin, C. Inducible expression
400 eliminates the fitness cost of vancomycin resistance in enterococci. *PNAS* **107**, 16964–
401 16969 (2010).

402

403 Correspondence and requests for materials should be directed to Kathy Raven.

404

405 **ACKNOWLEDGEMENTS**

406 We thank the Wellcome Trust Sanger Institute library construction, sequence and core
407 informatics teams, the staff at BSAC and the Cambridge Public Health England Microbiology
408 and Public Health Laboratory, and Hayley Brodrick, Amy Cain, Derek Pickard, Kim Judge and
409 Elizabeth Blane for their technical support. We thank the BSAC for allowing use of isolates
410 from the BSAC Resistance Surveillance Project. This publication presents independent
411 research supported by the Health Innovation Challenge Fund (HICF-T5-342 and WT098600), a
412 parallel funding partnership between the UK Department of Health and Wellcome Trust. The
413 views expressed in this publication are those of the authors and not necessarily those of the
414 Department of Health or the Wellcome Trust. This project was also funded by a grant
415 awarded to the Wellcome Trust Sanger Institute (098051). MET is a Clinical Scientist Fellow
416 supported by the Academy of Medical Sciences, the Health Foundation and the NIHR
417 Cambridge Biomedical Research Centre.

418

419 **AUTHOR CONTRIBUTIONS**

420 SJP designed the study. KER performed bacterial identification, susceptibility testing and DNA
421 extraction and analysed the data. SR assisted in bioinformatic analysis. TG, RR, JER, NMB and
422 JP contributed materials and data, MET completed ethical approvals, and JP and SJP were
423 responsible for supervision and management of the study.

424

425 **COMPETING FINANCIAL INTERESTS**

426 The authors declare that they have no competing financial interests.

427 **Figure legends**

428 **Figure 1. Phylogeny of *E. faecalis* isolates drawn from across the United Kingdom and**

429 **Ireland.** Left hand side: Midpoint rooted maximum likelihood tree of 168 *E. faecalis* isolates

430 based on SNPs in the core genome. Colored branches indicate the three dominant lineages

431 (L1, red; L2, purple; L3, turquoise). The vertical bars show the source of each isolate (dark

432 blue, BSAC; light blue, CUH; yellow, NCTC) and presence (red) or absence (blue) of

433 vancomycin resistance determinants. Bootstrap supports over 90% are labelled for the major

434 nodes. Scale bar indicates 10,000 SNPs. Right hand side: Maximum likelihood trees of the

435 three dominant lineages based on SNPs in the core genome after recombination was

436 removed, and rooted on an outlier. The trees are labelled by referral network, with ‘_1’ and

437 ‘_2’ indicating different hospitals within the referral network if more than one contributed to

438 the BSAC study collection, and year of isolation with CUH isolates highlighted blue. Bootstrap

439 supports over 90% are labelled. Scale bars indicate 25 SNPs.

440

441 **Figure 2: Global population structure of *E. faecalis*.** Phylogeny of 168 study isolates

442 combined with 347 isolates from geographically diverse locations downloaded from the

443 European Nucleotide Archive (ENA). Maximum likelihood tree based on SNPs in the 1,293

444 genes conserved in 99% of isolates. Colored branches indicate the three dominant lineages L1

445 (red), L2 (purple) and L3 (turquoise). Inner colored ring indicates the country of isolation,

446 outer colored ring indicates the source of the isolate.

447

448 **Figure 3: Prevalence of virulence and antibiotic resistance genes in the dominant lineages**

449 **(L1-3, n=89) and remainder (n=79).** Graphs show the percentage of isolates for which

450 putative virulence genes (a) or antibiotic resistance genes (grouped by antibiotic class) (b)

451 were detected. Genes that were ubiquitous in the collection are not shown. p values are

452 shown when a significant difference was observed using Fisher’s exact test. Virulence genes:

453 *ace* = collagen adhesion protein; *agg*= aggregation substance; *asa1*= aggregation substance;
454 *bgsA* = biofilm-associated glycolipid synthesis A; *cyl*= cytolysin; *elrA* = enterococcal leucine-
455 rich protein A; *esp*= enterococcal surface protein; *gelE* = gelatinase ; *perA* = pathogenicity
456 island-encoded regulator; *tpx* = thiol peroxidase. Antibiotic resistance genes: Am =
457 aminoglycosides (comprising one or more of *aac6'-2''*, *aph3''-III*, *aacA*, *ant-6-la*, *str*); Chlor =
458 chloramphenicol (*cat*); Linc = lincosamides (*InuB*); MLSB = macrolide, lincosamide,
459 streptogramin B (comprising *ermB* or *ermT*); Tet= tetracycline (comprising one or more of
460 *tetL*, *tetM*, *tetO*, *tetS*); Trim = trimethoprim (comprising *dfrC*, *dfrD*, *dfrF* or *dfrG*); Qac =
461 quaternary ammonium compounds and other antiseptics (*qacZ*), Vanc = vancomycin.

462

463 **Figure 4. Mapping variation in the vancomycin resistance transposon.** Midpoint rooted
464 maximum likelihood tree of all 168 *E. faecalis* with the 3 dominant lineages highlighted (L1,
465 red; L2, purple; L3, turquoise) and the presence (red) or absence (blue) of a *van* transposon
466 indicated in the vertical bar. The right hand side shows the coverage plot (number of
467 sequence reads that map to that location) of the *vanA* transposon in each isolate, with black
468 indicating presence (30x coverage or above), graduating to white indicating absence (less
469 than 10x coverage). The genes are labelled in the top bar (*tn* = inverted repeat). Colors in the
470 vertical bar on the right indicate the different insertion sites in the three dominant lineages,
471 with a description of each color available in Supplementary Table 4. Scale bar indicates
472 10,000 SNPs.

473