

1 **Evaluation of methods for whole genome sequencing of *Enterococcus faecium* in a**
2 **diagnostic laboratory**

3 **Running title: Laboratory methods for *E. faecium* sequencing**

4 Kathy E. Raven^{a*}, Danielle Leek^{a*}, Beth Blane^a, Sophia T. Girgis^a, Asha Akram^a, Nicholas
5 Brown^b, Sharon J. Peacock^a

6

7

8 ^aDepartment of Medicine, University of Cambridge, Box 157 Addenbrooke's Hospital, Hills
9 Road, Cambridge, CB2 0QQ, UK

10 ^bClinical Microbiology and Public Health Laboratory, Public Health England, Cambridge,
11 UK

12

13 *Joint first authors

14 Corresponding author: Kathy Raven (ker37@medschl.cam.ac.uk)

15

16 **Abstract**

17 *Enterococcus faecium* is an important nosocomial pathogen associated with hospital
18 transmission and outbreaks. Based on growing evidence that bacterial whole genome
19 sequencing enhances hospital outbreak investigation of other bacterial species, our aim was
20 to develop and evaluate methods for low volume clinical sequencing of *E. faecium*. Using a
21 test panel of 22 *E. faecium* isolates associated previously with hospital transmission, we
22 developed laboratory protocols for DNA extraction and library preparation, which in
23 combination with the Illumina MiniSeq can generate sequence data within 24 hours. The
24 final laboratory protocol took 3.5 hours and showed 98% reproducibility in producing
25 sufficient DNA for sequencing. Repeatability and reproducibility assays based on the
26 laboratory protocol and sequencing demonstrated 100% accuracy in assigning species,
27 sequence type (ST) and (when present) detecting *vanA* or *vanB*, with all isolates passing the
28 quality control metrics. Minor variation was detected in base calling of the same isolate
29 genome when tested repeatedly due to variations in mapping and base calling, but
30 application of a SNP cut-off (≤ 15 SNPs) to assign isolates to outbreak clusters showed 100%
31 reproducibility. An evaluation of contamination showed that controls and test *E. faecium*
32 sequence files contained <0.34% and <2.12% of fragments matching another species,
33 respectively. Deliberate contamination experiments confirmed that this was insufficient to
34 impact on data interpretation. Further work is required to develop informatic tools prior to
35 implementation into clinical practice.

36

37

38

39

40

41

42 **Importance**

43 *Enterococcus faecium* is a leading cause of hospital infections, particularly in the
44 immunocompromised, and has been shown to be associated with hospital transmissions.
45 Whole-genome sequencing is a highly discriminatory technique that has been shown to be
46 capable of identifying transmissions that may otherwise go undetected by conventional
47 infection control methods. This could be a powerful adjunct to infection control, since
48 previous studies have shown that there are numerous hospital outbreaks of *E. faecium*, which
49 can extend over multiple wards and years. In this study, we developed and evaluated
50 laboratory methods for low volume clinical sequencing of *E. faecium*, to allow uptake in
51 smaller local diagnostic laboratories. We demonstrated that sequencing could be performed
52 within 24-48 hours of the sample flagging positive. This would allow a more rapid
53 turnaround time, compared to sending isolates to the reference laboratory, providing an
54 opportunity for infection control to act earlier to prevent further transmission.

55

56 **Introduction**

57 *Enterococcus faecium* is a common gut commensal and a leading cause of nosocomial
58 infection, particularly in the immunocompromised (1). Infection caused by vancomycin-
59 resistant *E. faecium* (VRE) is of particular concern since there are limited treatment options
60 (2,3), as recognized by its inclusion in the WHO priority pathogens list for the research and
61 development of new antibiotics (4). The majority of nosocomial infection is associated with
62 hospital-adapted *E. faecium*, which are genetically distinct from commensal isolates (5).
63 Based on multilocus sequence typing (MLST), hospital-adapted strains cluster within Clonal
64 Complex (CC) 17, which has more recently been designated based on genome sequencing as
65 clade A1 (5). Isolates belonging to this clade are globally disseminated and have been linked
66 to numerous healthcare-associated outbreaks (6-8). Their enhanced biological fitness in
67 healthcare settings has been associated with genetic adaptations that alter survival, virulence
68 and antibiotic resistance (3).

69

70 The increasing application of whole genome sequencing to *E. faecium* has begun to elucidate
71 patterns of spread at national and individual hospital levels (6-8). This includes transmission
72 networks spanning multiple years and involving patient movement through numerous wards
73 within the same hospital, which are difficult for infection control to detect based on the
74 standard definition of same time and place (7). These studies have demonstrated that
75 bacterial sequencing can bring greater resolution to healthcare-associated lineages and
76 discriminate between isolates of the same lineage, overcoming this limitation of previous
77 typing methods such as pulsed-field gel electrophoresis (PFGE) and MLST. Due to
78 recombination involving regions containing the housekeeping genes used by MLST, this
79 method has been shown to be imperfect for determining genetic relatedness (8,9), and PFGE
80 has also been shown to cluster genetically distinct isolates (6). Our objective was to describe

81 the development and evaluation of laboratory methodology for low volume throughput
82 clinical sequencing of *E. faecium*.

83

84 **Materials and Methods**

85 *Test panel isolates*

86 Twenty-two *E. faecium* isolates were assembled into a test panel for the study (Table 1).
87 These were selected from a study published previously (7), and represented isolates
88 associated with bloodstream infection in patients at the Cambridge University Hospital NHS
89 Foundation Trust hospital (CUH) in the United Kingdom between Nov 2006 and Dec 2012.
90 Sixteen of these *E. faecium* were from four of six outbreaks described previously (7). The
91 remainder were two *E. faecium* belonging to each of Clade A2 (animal-associated clade) and
92 Clade B (community-associated clade); one *E. faecium* positive for *vanB*; and an *E. faecium*
93 positive control (Table 1). The collection included vancomycin-resistant and -susceptible
94 isolates.

95

96 *DNA extraction, library preparation and sequencing*

97 Isolates were cultured from frozen stocks onto CBA and incubated at 37°C for 48 hours. A
98 single colony was then sub-cultured onto CBA, incubated at 37°C for 48 hours and frozen at
99 -80°C in Microbank vials. For DNA extraction and subsequent sequencing, isolates were
100 taken from these pure stocks and incubated at 37°C on CBA overnight. DNA extraction was
101 performed manually using QIAgen kits according to the published QIAamp DNA Mini and
102 Blood Mini Handbook protocol, following ‘Appendix D: Protocols for Bacteria – Isolation
103 of genomic DNA from Gram-positive bacteria’ with the following amendments: (i) colonies
104 were used direct from the culture plate instead of pelleting bacteria, (ii) the lysozyme
105 solution was made using water and EDTA, (iii) the 95°C incubation was removed, (iv) in

106 steps 8 and 11 the centrifuge speed was set to maximum (13,200rpm), (v) the option of 50µl
107 of distilled water was used for elution including the 5 minute incubation step, with two 50µl
108 volumes filtered through in succession (final volume of 100µl), and 50µl of this used for a
109 final filter. Further amendments to the DNA extraction protocol were made as described in
110 the results. DNA was quantified post-extraction using a Qubit fluorometer and the dsDNA
111 HS assay kit (Thermofisher, UK). Library preparation was performed as described
112 previously using the Illumina Nextera Flex kit (10). Sequencing was performed on an
113 Illumina MiniSeq with a run-time of 13 hours, using the high output 150 cycle MiniSeq
114 cartridge and Generate Fastq workflow. Data were transferred to an external 1TB USB-
115 connected hard drive. Based on an expected total data output of 3.3-3.8Gb, an average
116 genome size of 2.9MB
117 ([https://www.ncbi.nlm.nih.gov/genome/?term=Enterococcus%20faecium\[Organism\]&cmd=](https://www.ncbi.nlm.nih.gov/genome/?term=Enterococcus%20faecium[Organism]&cmd=DetailsSearch)
118 [DetailsSearch](https://www.ncbi.nlm.nih.gov/genome/?term=Enterococcus%20faecium[Organism]&cmd=DetailsSearch)), and a target of 50x coverage, we estimated that approximately 24 *E. faecium*
119 isolates could be sequenced on a single run (estimated 47-55x coverage). We therefore
120 included 21 test *E. faecium* isolates and three controls (*E. coli*, *E. faecium* and no template,
121 see below) per sequence run for the initial reproducibility runs.

122

123 *Sequence data analysis*

124 Multilocus sequence types (ST) of the *E. faecium* isolates were identified using ARIBA
125 version 2.12.1 as described at [https://github.com/sanger-pathogens/ariba/wiki/MLST-](https://github.com/sanger-pathogens/ariba/wiki/MLST-calling-with-ARIBA)
126 [calling-with-ARIBA](https://github.com/sanger-pathogens/ariba/wiki/MLST-calling-with-ARIBA). Species were determined using Kraken version 1
127 (<https://ccb.jhu.edu/software/kraken/>) with the miniKraken database available at
128 https://ccb.jhu.edu/software/kraken/dl/minikraken_20171019_8GB.tgz. The presence of
129 *vanA* and *vanB* was determined using Ariba with *vanA* from M97297 (positions 6979-8010)
130 and *vanB* from Aus0004 (accession number CP003351, positions 2839377-2840405) as

131 reference genes. All isolates were mapped to the *E. faecium* strain Aus0004 (accession
132 number CP003351) using SMALT (<https://www.sanger.ac.uk/science/tools/smalt-0>) with
133 mapping and base calling performed as described previously (11), with the following
134 modifications: kmer size 13, step size 6. The depth and percentage coverage of the mapping
135 reference were determined using the script available at [https://github.com/sanger-](https://github.com/sanger-pathogens/vr-codebase/blob/master/lib/VertRes/Pipelines/Mapping.pm)
136 [pathogens/vr-codebase/blob/master/lib/VertRes/Pipelines/Mapping.pm](https://github.com/sanger-pathogens/vr-codebase/blob/master/lib/VertRes/Pipelines/Mapping.pm). Mobile genetic
137 elements (MGEs) were removed using the script available at [https://github.com/sanger-](https://github.com/sanger-pathogens/remove_blocks_from_aln)
138 [pathogens/remove_blocks_from_aln](https://github.com/sanger-pathogens/remove_blocks_from_aln). SNPs were identified based on the following
139 parameters: minimum number of reads matching the SNP = 4; minimum number of reads
140 matching the SNP per strand = 2; ratio of SNP base to alternative base >0.75; variant quality
141 >50; mapping quality >30.

142

143 SNP distances between isolate pairs reported previously (7) were defined after removal of
144 recombination regions using Gubbins across a large collection of clade A1 isolates. We
145 elected not to remove recombination in this study because looking to the future when
146 genomes are analysed using fully automated tools, recombination removal will currently be
147 challenging to incorporate into tools that offer very rapid interpretation. To correct for this
148 during the reproducibility of the pairwise SNP distance with the original sequence dataset,
149 we recalculated this based on mapping and removal of MGEs, but without regions of
150 recombination removed. These updated values were used as the ‘expected’ number of SNPs
151 between isolates for the repeatability and reproducibility runs.

152

153 ***Positive and negative controls***

154 Three controls were included in every sequencing run to monitor the ongoing performance
155 of the entire testing process. These were a no template control, a positive control (*E. faecium*

156 EC0160), and a negative control (*E. coli* NCTC12241). We selected a positive control that
157 was a known number of core genome SNPs (n=31) from another isolate sequence (EC0037)
158 from the same collection, to control for base calling in the place of PhiX. The positive
159 control was used to control the entire assay process and analytical accuracy. The negative
160 control was used to assess cross-contamination during processing and represented the non-
161 target DNA sample to verify analytical specificity. Fresh stocks of molecular grade water
162 and phosphate-buffered-saline were opened each week. The no template control contained
163 all assay components except for DNA and was used to verify the lack of contamination
164 across reagents and samples. Other ‘reuse’ reagents were checked for bacterial
165 contamination weekly by sub-culturing using a 10µl loop onto CBA and incubating for 48
166 hours in air at 37°C.

167

168 *Sequence metrics for controls*

169 Controls were required to pass the following quality metrics. *E. faecium* positive control:
170 highest match to *E. faecium* using Kraken, assigned to ST203, *vanA* detected, minimum
171 mean sequence depth of 20x and minimum 80% coverage of the mapping reference genome
172 (Aus0004). *E. coli* negative control: highest species match to *E. coli* in Kraken, *vanA* not
173 detected, no *E. faecium* ST assigned. No template control: contamination from any bacterial
174 DNA of less than 30,000 fragments in Kraken. *E. faecium* isolates from the test panel were
175 required to pass the following metrics: highest match to *E. faecium* using Kraken, assigned
176 to the correct ST, *vanA* detected or not detected as appropriate (Table 1), minimum sequence
177 depth of 20x and minimum 70% coverage of the mapping reference genome. This lower
178 value of 70% coverage was used for test isolates since Clade A2 and Clade B isolates are
179 more distantly related to the Clade A1 reference genome. Analysis of 799 genomes reported
180 previously (7,8) identified 17 genomes with <80% coverage of Aus0004 (74.8-79.6%

181 coverage), which belonged to Clade B (n=10), Clade A2 (n=4) or Bayesian Analysis of
182 Population Structure (BAPS) group 5 (n=3).

183

184 ***Repeatability and reproducibility***

185 Repeatability was evaluated by sequencing six *E. faecium* isolates (EC0160 (positive
186 control), EC0102, EC0181, EC0333, EC0397 and EC0503 (Table 1)) in triplicate in a single
187 sequencing run. Reproducibility was evaluated by sequencing 21 *E. faecium* isolates from
188 the test panel in three independent runs, and subsequently repeated with 12 *E. faecium*
189 isolates in three independent runs (see results). For each isolate, the pure single-colony
190 frozen stock was sub-cultured onto three separate CBA plates and incubated in air at 37°C
191 overnight, two heaped 1µl loops (Supplementary Figure 1) from each of these plates were
192 then taken forwards for individual DNA extraction, library preparation and sequencing. The
193 entire process for the reproducibility experiments was performed by different laboratory staff
194 on three different days. The resulting fastq files were analysed as above.

195

196 Isolates were classified as part of the same cluster if they were ≤ 15 SNPs apart. This cut-off
197 was selected based on a mutation rate of 7 SNPs/genome/year (7) and a within-host diversity
198 of 6 SNPs (12,13), using the formula described previously (14) to capture transmission
199 within 6 months. Isolates >15 SNPs different were classified as genetically unrelated.
200 Sensitivity and specificity for allocation of isolates into outbreaks were calculated using the
201 following definitions: true positives, the number of genetically related isolates based on the
202 original data that cluster together based on the test data; false negatives, the number of
203 genetically related isolates based on the original data that do not cluster together in the test
204 data; true negatives, the number of genetically unrelated isolates based on the original data

205 that do not cluster together in the test data; and false positives, the number of genetically
206 distant isolates based on the original data that cluster together based on the test data (15).

207

208 *Analysis of contamination*

209 The impact on quality metrics from varying levels of DNA contamination during clinical *E.*
210 *faecium* sequencing was evaluated using intentional spiking experiments. The *E. faecium*
211 positive control (EC0160), the *E. coli* negative control (NCTC 12241), and an *E. faecalis*
212 isolate (NCTC 13779) (selected because *E. faecalis* are commonly found on clinical plates
213 with *E. faecium*) were cultured and DNA extracted and quantified as described above. Donor
214 (contaminating) DNA was inoculated into the recipient (true) sample to achieve a final
215 spiked concentration of 0%, 0.1%, 1%, 10% or 20%. The donor-recipient combinations were
216 as follows: (i) recipient EC0160, donor NCTC 12241; (ii) recipient EC0160, donor NCTC
217 13779; (iii) recipient water, donor EC0160. Contamination with the spike was defined based
218 on the number and proportion of fragments matching to *E. faecium*, *E. faecalis* or *E. coli*
219 based on Kraken. The effect of contamination was evaluated using this metric together with
220 the proportion of the *E. faecium* reference covered during mapping, depth of coverage of the
221 mapping reference, and *vanA* and ST detected by Ariba. Additional unintentional
222 contamination from internal controls or external sources was evaluated based on the number
223 and proportion of reads matching to other species in Kraken.

224

225 *Data availability*

226 Sequence data generated during this study are available from the European Nucleotide
227 Archive (<https://www.ebi.ac.uk/ena>) under the accession numbers listed in Table 1.

228

229 **Results**

230 We sought to develop and describe methods for low-throughput *E. faecium* sequencing in a
231 routine microbiology laboratory within a 24-hour turnaround time (from DNA extraction to
232 availability of sequence data). This included an evaluation of quality controls, precision
233 (reproducibility and repeatability), and contamination.

234

235 First, we determined whether it was possible to extract DNA from a colony picked from the
236 primary clinical culture plate. Following the default DNA extraction protocol, we found that
237 extraction from a single colony after either 24 hours or 48 hours of incubation did not
238 provide enough DNA for input to the library preparation protocol, defined as being less than
239 the 3.3ng/μl required (Supplementary Table 1). We concluded that using colonies from the
240 primary plate was not feasible, and sub-cultured a single colony onto CBA and incubated for
241 a further 24 hours to create purity plates. DNA extraction from these purity plates
242 demonstrated that a single heaped 1μl loop input (Supplementary Figure 1) produced
243 insufficient DNA in 10% of cases (Supplementary Table 1), whilst two heaped 1μl loops
244 provided the required amount of DNA in all cases (Supplementary Table 1). We therefore
245 proceeded with two 1μl heaped loops input to DNA extraction.

246

247 We next aimed to determine whether the DNA extraction protocol time could be reduced
248 from the current time of 2 hours. Comparison of the DNA output using 30-minute
249 incubations for proteinase K and buffer AL versus 15-minute incubations at these steps
250 revealed that both methods resulted in sufficient DNA (Supplementary Table 1). We
251 therefore proceeded with a final protocol of two heaped 1μl loops input and 15 minutes
252 incubation for proteinase K and buffer AL, which reduced the time for DNA extraction from
253 2 to 1.5 hours. The final DNA extraction protocol was performed three times by three

254 different people, which demonstrated 98% reproducibility (65/66) for acquiring sufficient
255 DNA for input to library preparation (Supplementary Table 1).

256

257 Using the shortened DNA extraction protocol and previously described reduced library
258 preparation protocol (10), we aimed to determine the repeatability and reproducibility of the
259 full sequencing protocol. Repeatability was based on concordance of assay results and
260 quality metrics for six *E. faecium* isolates sequenced in triplicate in a single sequencing run.
261 There was 100% concordance in assigning species, ST and detecting *vanA* and *vanB*
262 (Supplementary Table 2). Analysis of the pairwise SNP differences between the within-run
263 replicates found that four of the six isolates were genetically identical across all three
264 replicates, whilst the remaining two isolates (EC0397 and EC0503) had 0-2 and 0-5 SNPs,
265 respectively, different between the three replicates. In the case of EC0397, replicates 1 and 3
266 were identical, but replicate 2 differed by 1-2 SNPs, whilst for EC0503 replicates 1 and 2
267 were identical but replicate 3 differed by 4-5 SNPs. This provided a repeatability per
268 replicate of 78% (14/18) based on a requirement for isolates to be identical, increasing to a
269 repeatability per replicate of 89% (16/18) if small variations in SNPs (≤ 2 SNPs) were
270 allowed. Analysis of the sequence files of EC0397 revealed that there were three positions
271 where the base calls varied between replicates. At two of these positions, both bases were
272 detected but in different proportions, while the third variable position could be explained by
273 misalignment around an indel. By contrast, EC0503 was identical in run 1 and run 2, but 5
274 SNPs different in run 3. Four of the five SNPs were located in a single 5bp region and were
275 likely caused by misalignment around an indel. Using the original published sequence
276 mapped to the Aus0004 reference with MGEs removed as the 'expected' number of SNPs,
277 the 6 isolates in triplicate had between 0-6 base calls different to the original sequence

278 (excluding positions denoted as ‘N’ because of failure to call a base), which is within the
279 expected within-patient diversity of 6 SNPs (Supplementary Table 2).

280

281 We initially evaluated reproducibility by sequencing 21 test panel *E. faecium* isolates in
282 three independent sequence runs (Supplementary Table 2). However, the average depth of
283 coverage across the three runs was lower than expected (43x, range 16.4-72.1x) despite a
284 higher data output than expected (4-5.2Gb compared to an expected 3.3-3.8Gb). Calculation
285 of the estimated genome size based on the data output and average depth of coverage across
286 the three sequence runs indicated a genome size of ~4.8Mb (Supplementary Table 3).
287 Investigation of the mapping files revealed that this could be caused by high depth of
288 coverage of plasmid elements, possibly due to multiple plasmid copies. Based on an
289 expected data requirement of ~4.8Mb, we estimated that 12 isolates plus three controls per
290 sequencing run would produce sufficient data to obtain 50x coverage. Repeats of the
291 reproducibility experiments with 12 isolates and 3 controls revealed an average depth of
292 coverage of 88.6x, with a range of 43.2-113.6x (Supplementary Table 2) based on 5.4-6.5Gb
293 data output, compared to an expected 89.9x coverage based on the observed data output.

294

295 Across the reproducibility runs with 12 isolates and 3 controls there was 100% accuracy in
296 assigning species, ST and detecting *vanA*. There were 0-3 SNPs identified for between-run
297 replicates, providing a reproducibility per replicate of 67% (26/39) based on replicates
298 needing to be identical, and a reproducibility per replicate of 92% (37/39) if small variations
299 (≤ 2 SNPs) were allowed between runs. Two isolates (EC0130 and EC0142) had 3 SNPs
300 different between runs. Analysis of the genomes revealed that in EC0130 these SNPs were
301 clustered in a 60bp region where a short section of reads had mapped in a region with poor
302 or no mapping, and both bases were present in the read files in all three repeats. In EC0142

303 two SNPs were adjacent in the genome in a short intergenic region lacking mapped reads in
304 the bam file, whilst the third was located in a short poorly mapped region in an otherwise
305 absent gene where all three repeats had a mixture of base calls at the position. Using the
306 original published sequence when mapped to the Aus0004 reference with MGEs removed as
307 the 'expected' number of SNPs, the new sequences were 0-4 SNPs different, which is within
308 the expected within-host diversity (6 SNPs).

309

310 We next sought to determine the sensitivity and specificity for outbreak detection in each of
311 the three reproducibility runs, using the genetic relatedness established previously as the
312 gold standard. The 12 isolates represented four distinct outbreaks encompassing four
313 different STs (ST203, ST375, ST117 and ST132) identified during 6 years of genomic
314 surveillance. All 12 isolates (36 isolate pairs) were classed in the same relatedness category
315 (0-15 SNPs, >15 SNPs) in each of the three sequence runs. This provides a sensitivity and
316 specificity for outbreak detection between runs of 100%. However, when the reproducibility
317 data was compared to the original sequence data, there were 2/12 discrepancies in the
318 classification of isolate pairs (one pair from Cluster A, one pair from Cluster D). Both pairs
319 were called as related by the original data (7 SNPs) and unrelated by the new data (16-24
320 SNPs). Analysis of the SNP locations revealed that the majority of the discrepant SNPs
321 (10/11 and 21/23 discrepant SNP locations in Cluster D and A, respectively) were due to a
322 single poorly mapped region which contained bacteriocin and plasmid genes and was located
323 between two ISEfm transposases, which had been called as unknown bases in the original
324 sequence. It was also found that 4 of the 7 SNPs in the original sequence data between the
325 isolate pair in Cluster A were clustered in a single 72bp region with poor mapping.

326

327 To determine the impact of DNA contamination (for example from cross-contamination
328 during processing) we performed deliberate contamination experiments. Details of the donor
329 and recipient DNA, the concentrations of spiked DNA and our findings are summarized in
330 Table 2. Contamination of the no template control with increasing concentrations of *E.*
331 *faecium* DNA did not lead to the control erroneously passing the QC metrics for *E. faecium*
332 until the final spiked concentration reached >10%. This indicates that contamination of the
333 no template control at 1% (which equated to 30,452 fragments matching *E. faecium* in
334 Kraken) can be tolerated. Contaminating the positive *E. faecium* control with increasing
335 concentrations of *E. coli* and *E. faecalis* DNA demonstrated that this could tolerate up to
336 10% contamination (which equated to 8.92 or 11.34% fragments, respectively, in Kraken)
337 before the *E. faecium* QC metrics were not achieved. Contamination with up to 20% *E. coli*
338 or *E. faecalis* did not appear to affect SNP calling (Supplementary Table 4).

339

340 We also evaluated unintentional contamination in the seven runs (excluding the deliberate
341 contamination assay). The maximum proportion of fragments matching another species was
342 0.34% for the controls and 2.12% for the test isolates, with the highest matches being those
343 of related species such as *Enterococcus durans* and *Lactococcus lactis*. Based on the number
344 of fragments in Kraken for the no template controls and the proportion of fragments in
345 Kraken for the remaining sequences, this demonstrates that all controls and test isolates had
346 levels of contamination below 1% (Supplementary Table 2).

347

348 **Discussion**

349 We aimed to develop manual methods for low-throughput sequencing of clinical *E. faecium*.
350 This would allow uptake of clinical sequencing in smaller local diagnostic laboratories that
351 lack the capacity for high-throughput sequencing. At present, isolates suspected to be part of

352 an outbreak are sent to the public health reference laboratory for further testing, with results
353 taking approximately two weeks to be returned (16). However, as the cost of sequencing and
354 sequencers reduces, the possibility of in-house sequencing at local clinical laboratories will
355 become more viable. This will allow a more rapid turnaround time to detection of outbreaks,
356 providing the opportunity for infection control to act earlier to prevent further transmission.
357 This is important since previous studies have shown that there are multiple outbreaks in a
358 hospital based on bloodstream infections alone, which represent only the tip of the iceberg
359 for transmission (7). Here we have shown that sequencing can be performed within 24-48
360 hours of a sample flagging as positive, providing a rapid tool to aid infection control.

361

362 In the clinical laboratory, two important factors for consideration are turnaround time and
363 cost. We aimed to identify the shortest turnaround-time from clinical plates to sequence data.
364 We found that single colonies produced insufficient DNA for sequencing, likely due to the
365 small colony size, meaning that purity plates are required. Purity plates may already be
366 available in the clinical laboratory at the time the sample flags as positive (for example,
367 purity plates used for disc testing), allowing a 24 hour turnaround time, whilst the remainder
368 will require a 48 hour turnaround time. This still represents a faster turnaround time than
369 testing at the reference laboratory. We were also able to reduce the hands-on processing time
370 from 4 hours to 3.5 hours for DNA extraction and library preparation. Since library
371 preparation does not require normalization (a tricky and potentially time-consuming step),
372 this represents a rapid and simple method for sequencing. To reduce the cost of sequencing
373 per sample we aimed to maximise the number of *E. faecium* isolates per sequence run. We
374 found that 21 isolates and 3 controls passed the quality control metrics in 95% of cases and
375 had a lower than expected coverage depth, potentially due to multiple copy number plasmids

376 increasing the data requirements for each isolate. We therefore proceeded with 12 isolates
377 and 3 controls.

378

379 Whilst the sequence data had 100% repeatability and reproducibility for assigning species,
380 ST and *van* genes, SNP detection was more variable. The majority of variation within and
381 between runs was minimal (0-1 SNPs) leading to a repeatability of 89% and reproducibility
382 of 92% if small variations (≤ 2 SNPs) were allowed, and the sensitivity and specificity for
383 outbreak detection between runs was 100%. However, some isolates had up to 5 SNPs
384 different within and between runs, and the sensitivity for outbreak detection dropped to 67%
385 when compared to the original sequence data. These discrepancies were explained by issues
386 with mapping and variant calling such as misalignment around indels and mapping of short
387 segments of the genome with high-density SNPs, including regions associated with mobile
388 genetic elements. Options for reducing these SNP errors could be to include recombination
389 detection, although this would be challenging to implement in rapid automated analysis
390 tools; improve the MGE file; use a core genome scheme such as that suggested by Coll et al.
391 (14), although this will need to be tailored to Clade A1 since community isolates are
392 considered genetically distant enough to classify as a different species (17); or alter the
393 mapping quality filters to reduce the impact of poorly mapped regions. The latter has been
394 used previously to improve detection of heterozygous sites in MRSA, where sites are only
395 considered if they are >50bp apart. Further work is needed to determine the optimum SNP
396 analysis methodology for this species.

397

398 Finally, we determined the level of contamination that could be tolerated and evaluated the
399 quality control metrics. All controls, test panel isolates and clinical isolates passed the
400 required quality control metrics. We found through a deliberate contamination experiment

401 that the data could be contaminated by up to 10% before quality control metrics were failed,
402 whilst all test isolates and controls showed <1% contamination.

403

404 Our findings indicate that the methods evaluated here can provide high quality sequence
405 data, and represents the first step towards being able to perform real-time clinical sequencing
406 of *E. faecium*. However, further work is required to develop data interpretation software that
407 can resolve issues relating to SNP calling in the *E. faecium* genome. This will be required
408 before rapid automated analysis tools can be developed for easy interpretation by clinical
409 staff.

410

411 **Acknowledgements**

412 This publication presents independent research supported by the Health Innovation
413 Challenge Fund (WT098600, HICF-T5-342), a parallel funding partnership between the
414 Department of Health and Wellcome Trust. The views expressed in this publication are those
415 of the author(s) and not necessarily those of the Department of Health or Wellcome Trust.

416

417 **Conflict of interest**

418 SJP is a consultant to Next Gen Diagnostics and Specific Technologies. Other authors have
419 no conflicts of interest.

420

421

422 **References**

- 423 1. Arias CA, Murray BE. The rise of the Enterococcus: beyond vancomycin resistance.
424 Nat Rev Microbiol 2012. 10(4): 266-278.
- 425 2. Werner G, Coque TM, Hammerum AM, Hope R, Hryniewicz W, Johnson A, Klare I,
426 Kristinsson KG, Leclercq R, Lester CH, Lillie M, Novais C, Olsson-Liljequist B,
427 Peixe LV, Sodowy E, Simonsen GS, Top J, Vuopio-Varkila J, Willems RJ, Witte W,
428 Woodford N. Emergence and spread of vancomycin resistance among enterococci in
429 Europe. Eurosurveillance 2008. 13(47), 19046.
- 430 3. Cattoir V, Leclercq R. Twenty-five years of shared life with vancomycin-resistant
431 enterococci: is it time to divorce? J Antimicrob Chemother. 2013. 68(4): 731-742.
- 432 4. World Health Organisation. Available online: [https://www.who.int/news/item/27-02-](https://www.who.int/news/item/27-02-2017-who-publishes-list-of-bacteria-for-which-new-antibiotics-are-urgently-needed)
433 [2017-who-publishes-list-of-bacteria-for-which-new-antibiotics-are-urgently-needed](https://www.who.int/news/item/27-02-2017-who-publishes-list-of-bacteria-for-which-new-antibiotics-are-urgently-needed)
434 (accessed on 11/06/2021).
- 435 5. Lebreton F, van Schaik W, McGuire AM, Godfrey P, Griggs A, Mazumdar V,
436 Corander J, Cheng L, Saif S, Young S, Zeng Q, Wortman J, Birren B, Willems RJL,
437 Earl AM, Gilmore MS. Emergence of epidemic multidrug-resistant *Enterococcus*
438 *faecium* from animal and commensal strains. mBio 2013. 20;4(4):e00534-13.
- 439 6. Pinholt M, Larnar-Svensson H, Littauer P, Moser CE, Pederson M, Lemming LE,
440 Ejlersen T, Sondergaard TS, Holzkecht BJ, Justesen US, Dzajic E, Olsen SS,
441 Nielsen JB, Worning P, Hammerum AM, Westh H, Jakobsen L. Multiple hospital
442 outbreaks of *vanA Enterococcus faecium* in Denmark, 2012-13, investigated by
443 WGS, MLST and PFGE. J Antimicrob Chemother 2015. 70(9): 2474-82.
- 444 7. Raven KE, Gouliouris T, Brodrick H, Coll F, Brown NM, Reynolds R, Reuter S,
445 Torok ME, Parkhill J, Peacock SJ. Complex routes of nosocomial vancomycin-

- 446 resistant *Enterococcus faecium* transmission revealed by genome sequencing. Clin
447 Infect Dis 2017. 64(7):886-893.
- 448 8. Raven KE, Reuter S, Reynolds R, Brodrick HJ, Russell JE, Torok ME, Parkhill J,
449 Peacock SJ. A decade of genomic history for healthcare-associated *Enterococcus*
450 *faecium* in the United Kingdom and Ireland. Genome Res 2016. 26(10): 1388-1396.
- 451 9. Howden BP, Holt KE, Lam MMC, Seemann T, Ballard S, Coombs GW, Tong SYC,
452 Grayson ML, Johnson PDR, Stinear TP. Genomic insights to control the emergence
453 of vancomycin-resistant enterococci. mBio. 2013. 4(4): e00412-13.
- 454 10. Raven KE, Blane B, Leek D, Churcher C, Kokko-Gonzales P, Pugazhendhi D, Fraser
455 L, Betley J, Parkhill J, Peacock SJ. Methodology for whole-genome sequencing of
456 methicillin-resistant *Staphylococcus aureus* isolates in a routine hospital
457 microbiology laboratory. J Clin Microbiol. 2019. 57(6):e00180-19.
- 458 11. Klemm EJ, Shakoor S, Page AJ, Qamar FN, Judge K, Saeed DK, Wong VK,
459 Dallman TJ, Nair S, Baker S, Shaheen G, Qureshi S, Yousafzai MT, Saleem MK,
460 Hasan Z, Dougan G, Hasan R. Emergence of an extensively drug-resistant
461 *Salmonella enterica* serovar typhi clone harboring a promiscuous plasmid encoding
462 resistance to fluoroquinolones and third-generation cephalosporins. mBio 2018. 9(1):
463 e00105-18.
- 464 12. Gouliouris T, Coll F, Ludden C, Blane B, Raven KE, Naydenova P, Crawley C,
465 Torok ME, Enoch DA, Brown NM, Harrison EM, Parkhill J, Peacock SJ.
466 Quantifying acquisition and transmission of *Enterococcus faecium* using genomic
467 surveillance. Nat Microbiol 2021. 6(1):103-111.
- 468 13. Brodrick HJ, Raven KE, Harrison EM, Blane B, Reuter S, Torok ME, Parkhill J,
469 Peacock SJ. Whole-genome sequencing reveals transmission of vancomycin-resistant
470 *Enterococcus faecium* in a healthcare network. Genome Med. 2016 8(1): 4.

- 471 14. Coll F, Raven KE, Knight GM, Blane B, Harrison EM, Leek D, Enoch DA, Brown
472 NM, Parkhill J, Peacock SJ. Definition of a genetic relatedness cutoff to exclude
473 recent transmission of methicillin-resistant *Staphylococcus aureus*: a genomic
474 epidemiology analysis. *Lancet Microbe* 2020. 1(8): e328-335.
- 475 15. Kozyreva VK, Truong CH, Greninger A, Crandall J, Mukhopadhyay RC, Chaturvedi
476 V. Validation and implementation of clinical laboratory improvements act-compliant
477 whole-genome sequencing in the public health microbiology laboratory. *J Clin*
478 *Microbiol* 2017. 55(8):2502–2520
- 479 16. National Infection Service. Bacteriology Reference Department User Manual.
480 Version 13, October 2020. Available online:
481 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachm](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/926234/BRDW0078.13_BRD_User_Manual.pdf)
482 [ent_data/file/926234/BRDW0078.13 BRD User Manual.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/926234/BRDW0078.13_BRD_User_Manual.pdf) (accessed on
483 17/06/2021)
- 484 17. Palmer KL, Godfrey P, Griggs A, Kos VN, Zucker J, Desjardins C, Cerqueira G,
485 Gevers D, Walker S, Wortman J, Feldgarden M, Haas B, Birren B, Gilmore MS.
486 Comparative genomics of enterococci: variation in *Enterococcus faecalis*, clade
487 structure in *E. faecium*, and defining characteristics of *E. gallinarum* and *E.*
488 *casseliflavus*. *mBio* 2012 3(1): e00318-11.

489 Table 1. Panel of isolates used in this study

ID	Accession number	Species	Vancomycin resistance	ST	Reason for inclusion	Repeatability	Reproducibility	Contamination
EC0397	ERR1820598	<i>E. faecium</i>	VSE	203	Cluster A	Included	Included	
EC0102	ERR1820597	<i>E. faecium</i>	<i>vanA</i>	203	Cluster A		Included	
EC0333	ERR370022	<i>E. faecium</i>	VSE	203	Cluster A			
EC0503	ERR1820599	<i>E. faecium</i>	<i>vanA</i>	203	Cluster A		Included	
EC0130	ERR375305	<i>E. faecium</i>	<i>vanA</i>	375	Cluster B		Included	
EC0518	ERR388722	<i>E. faecium</i>	<i>vanA</i>	375	Cluster B	Included	Included	
EC0185	ERR377443	<i>E. faecium</i>	<i>vanA</i>	375	Cluster B	Included		
EC0163	ERR375422	<i>E. faecium</i>	<i>vanA</i>	375	Cluster B	Included	Included	
EC0129	ERR370001	<i>E. faecium</i>	<i>vanA</i>	117	Cluster C		Included	
EC0175	ERR375395	<i>E. faecium</i>	<i>vanA</i>	117	Cluster C			
EC0142	ERR375299	<i>E. faecium</i>	<i>vanA</i>	117	Cluster C		Included	
EC0146	ERR375367	<i>E. faecium</i>	<i>vanA</i>	117	Cluster C		Included	
EC0031	ERR370024	<i>E. faecium</i>	<i>vanA</i>	132	Cluster D		Included	
EC0028	ERR377427	<i>E. faecium</i>	<i>vanA</i>	132	Cluster D		Included	
EC0036	ERR369956	<i>E. faecium</i>	<i>vanA</i>	132	Cluster D		Included	
EC0046	ERR375456	<i>E. faecium</i>	<i>vanA</i>	132	Cluster D			
EC0181	ERR375360	<i>E. faecium</i>	<i>vanB</i>	203	<i>vanB</i> positive	Included		
EC0333	ERR370022	<i>E. faecium</i>	VSE	203	Clade A2			
EC0322	ERR375339	<i>E. faecium</i>	VSE	Novel	Clade A2			
EC0214	ERR375555	<i>E. faecium</i>	<i>vanA</i>	203	Cluster A			
EC0392	ERR375541	<i>E. faecium</i>	VSE	Novel	Clade B			
EC0160	ERR377442	<i>E. faecium</i>	<i>vanA</i>	203	Positive control, Cluster A	Included	Included	Included
NCTC12241	ERR718772	<i>E. coli</i>	None	Not applicable	Negative control	Included	Included	Included
No template	Not applicable	None	None	None	No template control	Included	Included	Included

490

491 Cluster names A-D indicate individual outbreaks of *E. faecium* described previously (7).

492

493 Table 2. Results of deliberate contamination experiments

494

Question	Recipient	Donor	Concentration of contamination	Number/proportion of reads matching donor species	Proportion of reads matching <i>E. faecium</i>	Coverage of mapping reference	Average depth	MLST	<i>van</i>	QC metrics passed for positive control
Determine the effect of contaminating the no template control with increasing concentrations of <i>E. faecium</i> DNA	No template	<i>E. faecium</i>	0%	317	n/a	1.1	0.0	ND	None	No
			0.1%	1725	n/a	6.1	0.1	ND	None	No
			1%	30,452	n/a	60.0	1.1	ND	Partial match	No
			10%	42,015	n/a	86.9	11.4	203	<i>vanA</i>	No
			20%	47,621	n/a	87.4	20.6	203	<i>vanA</i>	Yes
Determine the effect of contaminating the <i>E. faecium</i> control with increasing concentrations of another organism	<i>E. faecium</i> control	<i>E. coli</i> control	0%	0.01%	85.95%	88.2	53.0	203	<i>vanA</i>	Yes
			0.1%	0.14%	85.73%	88.2	53.4	203	<i>vanA</i>	Yes
			1%	0.99%	84.05%	87.9	42.1	203	<i>vanA</i>	Yes
			10%	8.92%	72.29%	88.2	49.7	203	<i>vanA</i>	Yes
			20%	16.43%	61.12%	88.3	50.0	203	<i>vanA</i>	Yes
	<i>E. faecium</i> control	<i>E. faecalis</i>	0%	0.18%	85.86%	88.0	47.2	203	<i>vanA</i>	Yes
			0.1%	0.35%	85.75%	87.9	44.0	203	<i>vanA</i>	Yes
			1%	1.37%	84.75%	88.3	54.0	203	<i>vanA</i>	Yes
			10%	11.34%	74.84%	88.5	43.5	203*	<i>vanA</i>	No
			20%	21.92%	64.12%	88.5	37.3	203*	<i>vanA</i>	No

495

496 * indicates that there was uncertainty in the MLST call by Ariba.

497

498 **Supplementary figure 1:** Figure showing a heaped 1µl loop of *E. faecium* used as input to DNA extraction

499

