1  **Accuracy of different bioinformatics methods in detecting antibiotic resistance**

2  **and virulence factors from *Staphylococcus aureus* whole genome sequences.**

3

4  Authors: Amy Mason*[1], Dona Foster*[1]#, Phelim Bradley*[2], Tanya Golubchik*[1], Michel

5  Doumith*[3], N. Claire Gordon[1], Bruno Pichon[3], Zamin Iqbal[2], Peter Staves[3], Derrick

6  Crook[1,4,5,6], A. Sarah Walker**[1,5,6], Angela Kearns**[3,5], Tim Peto**[1,5,6]

7  */** contribution considered equal

8

9  [1] Nuffield Department of Medicine, University of Oxford, UK

10  [2] Wellcome Trust Centre for Human Genetics, University of Oxford, UK

11  [3] Staphylococcus Reference Service, National Infection Service, Public Health England,

12  UK

13  [4] National Infection Service, Public Health England, UK

14  [5] The National Institute for Health Research (NIHR) Health Protection Research Unit in

15  Healthcare Associated Infections and Antimicrobial Resistance at University of Oxford,

16  UK

17  [6] NIHR Oxford Biomedical Research Centre, University of Oxford, UK

18

19  Running Head: *S. aureus* Whole-genome Sequence Method Comparison

20

21  #Address correspondence to Dr Dona Foster, Microbiology Level 7, John Radcliffe

22  Hospital, Headley Way, Oxford, OX3 9DU. dona.foster@ndm.ox.ac.uk

23  Current institution: Amy Mason: Department of Mathematics and Department of

24  Statistics, University of Oxford, UK. N. Claire Gordon: KEMRI-Wellcome Trust

25    Collaborative Research Programme, Kilifi, Kenya. Tanya Golubchik: Wellcome Trust

26    Centre for Human Genetics, University of Oxford, UK.

27

28    Length: <mark>2999</mark> words (limit 3000 excluding Materials and Methods), 3 Tables (plus 4

29    Supplementary), 3 Figures (plus 2 Supplementary)

30 **Abstract (249 words, limit 250 words)**

31 **Background**: In principle, whole genome sequencing (WGS) can predict phenotypic

32 resistance directly from genotype, replacing laboratory-based tests. However, the

33 contribution of different bioinformatics methods to genotype-phenotype discrepancies

34 has not been systematically explored to date.

35 **Methods**: We compared three WGS-based bioinformatics methods (Genefinder (read-

36 based), Mykrobe (de Bruijn graph-based) and Typewriter (BLAST-based)) for predicting

37 presence/absence of 83 different resistance determinants and virulence genes, and

38 overall antimicrobial susceptibility, in 1379 *Staphylococcus aureus* isolates previously

39 characterised by standard laboratory methods (disc diffusion, broth and/or agar

40 dilution and PCR).

41 **Results**: 99.5% (113830/114457) of individual resistance-determinant/virulence gene

42 predictions were identical between all three methods, with only 627 (0.5%) discordant

43 predictions, demonstrating high overall agreement (Fliess-Kappa=0.98, p<0.0001).

44 Discrepancies when identified were in only one of the three methods for all genes except

45 the cassette recombinase, *ccrC(b)*. Genotypic antimicrobial susceptibility prediction

46 matched laboratory phenotype in 98.3% (14224/14464) cases (2720 (18.8%) resistant,

47 11504 (79.5%) susceptible). There was greater disagreement between the laboratory

48 phenotypes and the combined genotypic predictions (97 (0.7%) phenotypically-

49 susceptible but all bioinformatic methods reported resistance; 89 (0.6%)

50 phenotypically-resistant, but all bioinformatics methods reported susceptible) than

51 within the three bioinformatics methods (54 (0.4%) cases, 16 phenotypically-resistant,

52 38 phenotypically-susceptible). However, in 36/54 (67%), the consensus genotype

53 matched the laboratory phenotype.

54  **Conclusions**: In this study, the choice between these three specific bioinformatic

55  methods to identify resistance-determinants or other genes in *S. aureus* did not prove

56  critical, with all demonstrating high concordance with each other and

57  phenotypic/molecular methods. However, each has some limitations and therefore

58  consensus methods provide some assurance.

59

**Introduction**

*Staphylococcus aureus* causes both superficial infections (such as boils) and life-threatening disease including septicaemia (1). There were 11,405 *S. aureus* bacteraemias in England in 2015/2016 (2); 7.2% were meticillin resistant *S. aureus* (MRSA) which has increased costs and poorer patient outcomes (3). Fast accurate resistance prediction is key to managing *S. aureus* infections. Molecular-based methods directed at detecting specific genes, e.g. through rapid multiplex PCR and microarrays, can reduce time to identify resistance determinants and time on broad-spectrum antibiotics (4-6). However, they require specific primers that impact sensitivity and specificity.

In principle, whole genome sequencing (WGS) has the potential to predict phenotypic resistance directly from genotype, replacing laboratory-based phenotypic tests (7). Several studies report high concordance between genotypic predictions based on known or novel resistant determinants and phenotypic methods (8-13). However, these studies used varying sequence processing pipelines and bioinformatics methods to identify *in silico* resistance determinants. Without formal comparisons between the various methods, it is unclear whether the underlying differences affect results, or whether differences in methodology could cause some of the observed discrepancies between genotypic predictions and phenotype.

Here, we therefore compare three WGS-based bioinformatics methods (Genefinder (read-based), Mykrobe (de Bruijn graph-based) and Typewriter (BLAST-based)) in terms of predictions of presence/absence of different resistance determinants, and

84    overall prediction of antimicrobial susceptibility and presence/absence of virulence

85    genes, from short-read Illumina WGS.

86

87    **Results**

88    Short-read Illumina WGS were available from 1,389 samples, 992 from a collection held

89    in Oxford (previously described by Gordon *et al* (9, 10)) and 397 from Public Health

90    England (PHE) Staphylococcus Reference Service, Colindale. Ten samples were excluded

91    due to mixed/contaminated WGS results, leaving 1,379 for analysis. Samples were

92    analysed by Genefinder and Typewriter (Table 1) after sequence mapping and variant

93    calling and by Mykrobe from raw fastq reads.

94

95    84 genes were included: 46 acquired resistance genes, five sets of chromosomal variants

96    within genes associated with resistance, three cassette chromosome recombinases *ccrA*,

97    *ccrB* and *ccrC* including three variants of *ccrC* (*ccrCa, ccrCb, ccrCc*) and 28 virulence

98    genes (Supplementary Table 1). 99.5% (113830/114457) of the individual resistance-

99    determinant/virulence gene predictions were identical between all three methods

100   (Supplementary Table 1, Figure 1), with only 627 (0.5%) discordant predictions,

101   demonstrating high overall agreement (Fliess-Kappa=0.98, p<0.0001). Overall, one

102   method disagreed with both other methods in 0.23% for Typewriter (263/114457

103   predictions), 0.16% Mykrobe (183/114457) and 0.16% Genefinder (181/114457). The

104   three most common discrepancies for Typewriter were the non-detection of virulence

105   genes identified by other methods (*seu* 57 samples, *chp* 46 samples, *sei* 33 samples).

106   Similarly, for Genefinder the three most common discrepancies were non-detection of

107   resistance genes (*qacB* 44 samples, *dfrC* 34 samples) or other genes (*ccBb* 22 samples)

108   identified by other methods. Genefinder reported the presence of *dfrA*, *qacA* or *ccrC*(b)

109    genes in these samples. In contrast, Typewriter and Mykrobe reported the presence of

110    two *dfr*, two *qac* and three *ccrC* genes, where the detected variants for each of these

111    three genes shared more than 90% nucleotide identity. The most common discrepancies

112    for Mykrobe were identifying resistance/other genes as present when the other two

113    methods called them absent (*aadE/ant(6)-Ia* 28 samples, *blaZ* 19 samples, *ccrCB* 22

114    samples). No gene was ever identified as present by Typewriter alone. 14 of the 84

115    genes had >1% discrepancies (maximum 4.3% for *seu*), but the majority of discrepancies

116    were in only one method for all genes except *ccrC(b)*.

117

118    Discrepancies were similar in acquired resistance genes (0.3%, 221/63434) and

119    chromosomal resistance genes (0.1%, 8/5516), but slightly larger for *ccr* genes (1.8%,

120    123/6895) and virulence genes (0.7%, 275/38612) (Supplementary Table 2).

121    Percentage discrepancies varied modestly across the different sample sets, being higher

122    for the PHE set (1.1%, 349/32,928; particularly for *ccr* genes with 4.2% (83/1,960)

123    discrepancies), intermediate for the Oxford derivation set (0.6%, 233/42084) and

124    lowest for the Oxford validation set (0.1%, 45/40,824) (Supplementary Table 2).

125

126    Genotypic predictions of antimicrobial susceptibility were also identical in 99.6% of

127    cases (16,477/16,548 predictions, Table 2). Of the 71 discrepancies in susceptibility

128    prediction between the methods, 42% (30/71) occurred with Typewriter reporting

129    susceptible when Genefinder and Mykrobe reported resistant, and 49% (35/71)

130    occurred with Mykrobe reporting resistant where Genefinder and Typewriter reported

131    susceptible.

132

133     Comparing genetic predictions to laboratory phenotypes (restricted to samples either

134     phenotypically resistant or susceptible), in 98.3% (14224/14464) cases all three

135     bioinformatics methods and the gold standard laboratory results agreed completely

136     (2720 (18.8%) resistant, 11504 (79.5%) susceptible) (Table 3a, Figure 2). There was

137     greater disagreement between the laboratory phenotypic results and the combined

138     genotypic predictions than within the three bioinformatics methods. In 97 (0.7%)

139     instances, the laboratory phenotype was susceptible but all bioinformatic methods

140     reported resistance. Of these, 33% (32/97) were for penicillin, 23% (22/97)

141     clindamycin and 11% (11/97) erythromycin, with smaller numbers for fusidic acid (7),

142     tetracycline (6), mupirocin (6), methicillin (5), ciprofloxacin (4), gentamicin (3) and

143     rifampicin (1), and none for trimethoprim. In 89 (0.6%) instances, the laboratory

144     phenotype was resistant, but all three bioinformatics methods reported susceptible,

145     most commonly to gentamicin (21%, 15/89), ciprofloxacin (17%, 15/89) and fusidic

146     acid (15%, 13/89). The remaining 54 (0.4%) cases (16 phenotypically-resistant, 38

147     phenotypically-susceptible) had different genotypic predictions made from the different

148     methods. However, in 36/54 (67%), the consensus genotype (predicted by two of the

149     three methods) matched the laboratory phenotype.

150

151     PCR/array results were available for some virulence genes (14) and *mecA/mecC* for all

152     397 PHE isolates. Compared with genetic predictions, in 96.8% (3983/4115) cases all

153     three bioinformatics methods and the PCR/array results agreed completely (3364

154     (81.7%) absent, 619 (15.0%) present) (Table 3b, Supplementary Figure 1). As for

155     antimicrobial resistance, there was greater disagreement between the laboratory

156     PCR/array results and the combined genotypic predictions than within the three

157     bioinformatics methods, with 81 (2.0%) cases where all three methods called a gene

158    present that had not been detected by PCR/array and 12 (0.3%) where no method called

159    a gene present that had been detected by PCR/array, in comparison with 39 (0.9%)

160    discrepant predictions between the methods. In 20/39 (51%), the consensus genotype

161    matched the PCR/array result.

162

163    The sensitivity and specificity of all three bioinformatics methods compared to

164    laboratory phenotypic methods in predicting antimicrobial susceptibility was very

165    similar. Across the 14464 genotypic predictions, Typewriter had the lowest overall

166    sensitivity (0.964 (95% CI 0.956-0.970), but the highest specificity (0.992 (0.990-

167    0.993)), while Mykrobe had higher sensitivity (0.967 (0.960-0.974)) and lowest

168    specificity (0.989 (0.987-0.990)). Genefinder's performance fell between Mykrobe and

169    Typewriter for specificity (0.990 (0.988-0.992)) with a sensitivity equal to Mykrobe

170    (0.967 (0.960-0.973)). Specificity and sensitivity varied across the different antibiotics

171    (Figure 3), but were broadly similar between the three methods, overall and within each

172    dataset (Supplementary Table 3). There were no vancomycin resistant isolates

173    identified by either phenotyping or bioinformatics methods. Similarly, specificity and

174    sensitivity to identify PCR-detected virulence and other genes varied across the different

175    genes, but were broadly similar between the three methods (Supplementary Figure 2).

176

177    **Discussion**

178    Whilst WGS is increasingly used to detect antibiotic resistance and virulence

179    determinants, to our knowledge this is the first study that compares three methods for

180    predicting genotype on large numbers of isolates. As discussed in the recent European

181    Committee on Antimicrobial Susceptibility Testing (EUCAST) report (15), discordance

182    can occur between phenotypic and genotypic resistance due to inadequate limits of

183     detection for WGS methods, incomplete understanding of the genotypic basis of

184     phenotypic resistance, flaws with the phenotypic or molecular (e.g. PCR) methods

185     currently used to detect resistance, and/or WGS failures including lack of assembly

186     caused by multiple operons or similar sequences, incomplete gene coverage, non-

187     functional genes (e.g., due to presence of stop codons/indels) or cropped contigs.

188

189     Here we found that three different approaches to identifying genetic determinants of

190     resistance and virulence (Genefinder, Mykrobe and Typewriter) agreed in 99.5%

191     predictions. Genefinder and Mykrobe were fast, taking under five minutes whereas

192     Typewriter, while also taking a few minutes per sample, required initial genome

193     assembly that increased turnaround time by up to three hours. Mykrobe and Typewriter

194     are freely available (https://github.com/iqbal-lab/Mykrobe-predictor and

195     https://github.com/tgolubch/typewriter respectively); Genefinder is not but the

196     underpinning methods are relatively straightforward, and the freely available SRST2

197     (https://github.com/katholt/srst2) follows an analogous mapping approach (16) which

198     would likely provide very similar results with the same catalogue. Previous comparisons

199     of bioinformatics methods relevant to the microbiology community are limited. Bradley

200     et al (9) found good concordance between Mykrobe and SeqSphere (17), an allele-based

201     method that detects presence/absence of a limited number of resistance and virulence

202     markers. SeqSphere took longer than Mykrobe as, like Typewriter, it uses Velvet

203     assemblies. Other previous studies have shown 100% concordance between resistome

204     and toxome in 14 MRSA isolates (18), 98.6% concordance across 5288 susceptibility

205     predictions in 308 *S. aureus* isolates (both MRSA and MSSA) (19), 100% concordance for

206     selected resistance and toxin gene presence/absence in 18 MRSA strains (17), and

207     97%/97% sensitivity/specificity for Typewriter and 99.1%/99.6%

208    sensitivity/specificity for Mykrobe for predicting phenotypic resistance in the Oxford

209    validation samples used here (9, 10). A comparison between microarray and WGS in 154

210    isolates reported 1.7% discordancy in detecting resistance and virulence genes (20),

211    mainly due to failure of WGS to detect enterotoxins and super antigens (similar to

212    Typewriter in this study).

213

214    Individually, the three programs demonstrated high concordance, but interestingly, in

215    almost all genes only one of the three bioinformatics methods did not identify a

216    determinant that the other two methods did identify, or vice versa. The most common

217    discrepancy with Typewriter was failing to identify virulence genes identified by

218    Mykrobe and Genefinder (namely, *seu*, *chp* and *sei*). Two of these genes, *sei* and *seu*, are

219    located on the enterotoxin gene cluster (*egc*) (21, 22), referred to as an enterotoxin gene

220    nursery (23), and the other, *chp*, on a prophage (24). Such regions may be particularly

221    susceptible to recombination (25, 26) and paralogs. As Typewriter uses BLAST, it may

222    have a higher chance of detecting one of multiple closely related genes than the other

223    two methods.

224

225    Similarly to Typewriter, the most common discrepancy with Genefinder was failing to

226    identify genes reported by Typewriter or Mykrobe, particularly *ccrB, qacB* (*quaternary*

227    *ammonium compound B*, conferring resistance to chlorexidine (27) via an efflux drug

228    pump, but differing from another gene, *qacA*, by only seven nucleotides (28)), and *dfrC*

229    (a dihydrofolate conferring resistance to trimethoprim believed to be the origin of the

230    more common transposon-associated *drfA* gene). The fact that Genefinder identified

231    only one variant of acquired *dfr* and *qac* may indicate that the other two methods were

232    misidentifying paralogs (29). Alternatively, as Genefinder detects pre-determined

233    alleles, recombination of partial genes or differences in flanking sites or genomic

234    variation alone may reduce its ability to detect some genes. One advantage of Genefinder

235    is its ability to detect variations in multicopy genes such as the ribosomal RNA encoding

236    genes associated with linezolid resistance in staphylococci.

237

238    In contrast, Mykrobe most commonly identified a determinant that other methods did

239    not, particularly *aadE(ant6')-Ia,* an adenyltransferase encoding resistance to

240    aminoglycosides. This gene is associated with small plasmids flanked by direct repeats

241    of staphylococcal insertion sequence IS257 (30). Although Mykrobe is kmer-based, it

242    requires a high match across the whole gene, not just flanking sequences, so the reason

243    for this is unclear. Mykrobe also had a higher false-positive rate in *blaZ*, as reported

244    previously (9). Although this was previously attributed to phenotypic errors, the fact

245    that neither Genefinder nor Typewriter identified *blaZ* in these isolates suggests the

246    algorithm/threshold may need adjusting for this gene. Mykrobe also had a high false-

247    positive rate for the *ccrCB* gene, which is part of the cassette chromosome recombinase

248    (*ccr*) associated with SSC*mec* (31). As all *ccrC* genes share >87% similarity, and were not

249    included in the original Mykrobe implementation, further investigation and modification

250    of sequence identity thresholds may be required to accurately classify this gene, whose

251    different alleles can have 60-82% sequence identity.

252

253    Overall, the comparison highlights key challenges inherent in all methods. First is the

254    trade-off between specificity and sensitivity to detect specific genes/variants, and the

255    need for adjustment based on specific features, such as proximity to repetitive elements

256    or similarity with other alleles. Specific genes may also require different approaches, e.g.

257    the *ccr* genes were the most discordant overall in the study. These genes were more

258    often present in the Staphylococcal reference laboratory isolates, increasing overall

259    error rates for this sample set. Reference libraries of genes/variants also require

260    frequent updating with new alleles, and appropriate thresholds must be set to allow

261    separate copies of closely related genes (e.g. *qacA* and *qacB*) to be detected if genuinely

262    present. Taking the consensus prediction across the three different bioinformatics

263    methods is one strategy for balancing these different trade-offs. As error rates were low

264    overall, this only improved genetic predictions slightly, but in samples where the

265    susceptibility is unknown it could be valuable, particularly if the two fast

266    implementations (GeneFinder, Mykrobe) are used, followed by the slower assembly-

267    based method only if they disagree.

268

269    Our main findings were that the largest discordance occurred between phenotype and

270    genotype regardless of the method used to predict genotype, and that the "consensus"

271    genotypic prediction agreed with the phenotype in two-thirds of the small number of

272    cases where bioinformatics methods made different predictions. Where bioinformatics

273    methods are concordant, but disagree with phenotype, the unresolved question is which

274    is "correct", in terms of a drug achieving clinical cure in a patient infected with this

275    strain. Penicillin and clindamycin/erythromycin were most likely to be called resistant

276    by all methods but susceptible by phenotyping. Previous studies of erythromycin and

277    clindamycin resistance have reported positive *ermC* PCR results from non-detectable

278    resistance phenotypes (32) and have suggested that plasmids conferring resistance to

279    these antibiotics may be lost in subculture (9, 33). Sensitivity to penicillin by phenotypic

280    methods where genotype methods predict resistance has been reported previously (34,

281    35) and the evidence suggests that phenotyping underreports resistance. The EUCAST

282    guidelines illustrate the challenges in distinguishing between penicillin-resistant and -

283    susceptible isolates based on fuzzy versus sharp zones (36). Overall therefore it is

284    plausible that genetic detection of resistance may reflect more closely the impact of the

285    strain on a patient.

286

287    Interpretation where phenotyping reports resistance but WGS methods predict

288    susceptibility is more difficult. One possibility is small colony variants (SCV) being

289    present phenotypically but overgrown in WGS culture and thus not represented in the

290    sequence. Resistance associated with gentamicin, fusidic acid and ciprofloxin, the main

291    antibiotics where this phenomenon was observed, is observed with SCV phenotypes (37,

292    38). An alternative explanation is novel resistance mechanisms, for example,

293    ciprofloxacin (39), leading to false-negative WGS predictions. The need for a

294    continuously updated curated database is a key challenge for WGS methods. As more

295    sequencing occurs, novel mutations will be identified in resistance genes that may or

296    may not confer phenotypic resistance, but these can at least be identified and tested;

297    identifying entirely new resistance-conferring genes is more complex and prediction

298    software that can recognize new, clinically important genes a priori would be a valuable

299    addition to an analysis pipeline. However, we observed similar differences between

300    concordant genotypic predictions and both phenotypic antimicrobial susceptibilities

301    and single gene PCR results, suggesting that the underlying causes may not necessarily

302    be related to resistance. As previously noted, agreement between WGS and phenotyping

303    is higher (98.6%) than between phenotyping undertaken by two separate laboratories

304    (97.6%) (19), thus at least some discrepancies are probably due to incorrect

305    phenotyping results. In contrast, concordance between genotypic predictions made

306    using a single method but based on WGS generated from 5 different laboratories was

307    recently shown to be >99.8% (40).

## Limitations

This comparison was based on a pre-specified set of resistance or virulence associated genes: some genetic traits previously associated with resistance were omitted (eg. *IleS* mutations linked to low-level mupirocin resistance). Despite this, we found good agreement between genotypic predictions and phenotype. Typewriter used Velvet de novo assemblies: other newer assemblers (e.g. SPADES (41)) might have improved predictions further. We included data which had been used in development of two of the methods compared, which could potentially have led to over-fitting, although performance of all three methods was in fact similar on this dataset (Supplementary Table 3). All analysis was undertaken on short-read Illumina data. The increasing use of long-read sequences will require further software testing, although Mykrobe has been successfully used for initial resistance calling in *Mycobacterium tuberculosis* from Nanopore sequencing in a small number of samples (42). However, it has not been comprehensively tested, nor have Typewriter or Genefinder, with long-read sequences generated using Nanopore or PacBio technology. The greatest differences detected in this study were between phenotype and genotype, which could be partly due to the method of phenotypic testing and recognised issues with reproducibility. We did not have resources to re-phenotype all or a subset of the isolates; well-characterised sets of repeatedly phenotyped isolates would be useful for further studies. We found no suggestion that missing calls in one program were associated with scores just below a threshold, but did not undertake a more detailed assessment of specific sequence coverage and quality around discrepant genetic predictions.

## Conclusion

331 In summary, in this study the choice between three specific bioinformatic methods to

333 identify resistance-determinants or other genes in *S. aureus* did not prove critical. All

334 demonstrated a high concordance with each other, and phenotypic methods, and can be

335 recommended for genotype prediction. However, each has some limitations and

336 therefore consensus methods provide at least some assurance. Due to computational

337 speed, Mykrobe (de Bruijn graph-based) and Genefinder (or equivalent mapping-based

338 program such as SRST2 (16)) are a sensible combination to use as an initial consensus

339 method, followed by Typewriter (BLAST-based) if these two methods disagree. As a set

340 of 34 diverse bacteria have been made available for whole genome sequencing

341 validation (43), the study strains and genotypic predictions are available as a resource

342 for other studies investigating different bioinformatic analysis methods which will

343 become increasingly important as this technique is more widely used to inform clinical

344 management, though bacterial identification, antimicrobial susceptibility prediction and

345 virulence profiling. External quality control of clinical laboratory performance in

346 predicting antibiotic resistance is provided by UK proficiency testing schemes such as

347 UK NEQAS (United Kingdom National External Quality Assessment Service for

348 Microbiology) (44); a similar set of standards will need to be created to accredit whole

349 genome sequencing methods.

350

## Materials and Methods

352 Three sets of *S. aureus* isolates with known high-quality phenotypes were analysed: a

353 derivation, n=501, and validation, n=491, set (denoted "Oxford derivation/validation")

354 from blood cultures and nasal swabs isolates at the Oxford Radcliffe Hospitals NHS Trust

355 and Brighton and Sussex University Hospitals NHS Trust, spanning a period of 13 years,

356    sequenced for an initial assessment of genotypic prediction of susceptibility phenotype

357    in *S. aureus* (9, 10)and 397 isolates that had been referred to the Public Health England

358    reference laboratory for investigation (denoted "Colindale 397", available at NCBI:

359    PRJNA445516). The Oxford derivation set had previously been used in the development

360    of Typewriter and Mykrobe, but not Genefinder; the former methods were then applied

361    to the Oxford validation set.

362

363    Phenotypes for "Oxford derivation/validation" isolates used disc diffusion and/or

364    automated broth diffusion (BD Phoenix) with discrepancies between phenotype and

365    genotype resolved as described previously (11). All PHE isolates (n=397) were

366    subjected to MIC testing by the PHE Staphylococcal Reference Laboratory using the agar

367    dilution method (45). In addition, the *mec*A/C status and virulence gene profile of the

368    PHE isolates was determined by PCR or microarray testing as described previously (14).

369    The European Committee on Antimicrobial Susceptibility Testing (EUCAST): thresholds

370    were used to determine sensitivity or resistances for each phenotype

371    (http://www.eucast.org/clinical_breakpoints).

372

373    All "Oxford derivation/validation" isolates were sequenced using the Illumina HiSeq

374    2000 platform as previously described (46). PHE samples were sequenced in an

375    Illumina HiSeq 2500 platform as described previously (47) (both 150bp reads). Samples

376    determined as mixed based on WGS were excluded from further analysis. Quality control

377    of sequences at PHE used the trimmomatic software (Illumina adapter removed, leading

378    and trailing quality threshold set to 30 and minimum length of read set to 50 bases)

379    (48).  Isolates from Oxford analysed by Typewriter were mapped and de novo

380    assembled with exclusion parameters of <70% coverage of reference genome for

381  mapping and <50% of the genome in contigs >1 Kb (10). Mykrobe processes raw

382  sequence data with no prior cleaning of the data. Isolates came from 111 sequence

383  types, including 29 new STs/alleles, covering the range of *S. aureus* genomic diversity as

384  previously described in Oxfordshire.

385

386  Three programs, Genefinder (MD; PHE, not published), Mykrobe (PB; Version v0.3.13-2-

387  gd5880fa, open-source at https://github.com/iqbal-lab/Mykrobe-predictor), and

388  Typewriter (TG; version 2.0, MMM group, Oxford University,

389  https://github.com/tgolubch/typewriter) (Table 1), were compared to determine

390  presence/absence of resistance-determinants (genes or variants) and toxin genes

391  (Tables 2, 3). Mykrobe is part of the automated processing with the Complete Pathogen

392  Software Solution (COMPASS) developed at University of Oxford. This returns quality

393  and depth of sequence metrics, maps against a reference (MRSA 252, GenBank

394  Accession no: BX57186561) using Stampy (49) and performs *de novo* assembly using

395  Velvet v1.0.18 (50). These de novo assemblies formed the basis for the Typewriter

396  program, whereas Genefinder used the raw sequencing reads.

397

398  Although all three methods search for matches to a pre-defined list of alleles, they have

399  different approaches to their identification (further details below). Genefinder and

400  Mykrobe required fastq files whereas Typewriter used BLAST on de novo assemblies. All

401  used pre-set thresholds to detect genes. Thresholds are adapted for certain genes (e.g.

402  *blaZ* which can be chromosomally integrated or carried on plasmids) to improve

403  prediction and for quality control. Both Typewriter and Mykrobe identified presence or

404  absence of each target singly, whereas Genefinder identified which of closely related

405  homologs is most plausibly present. Genefinder and Mykrobe were very fast, between

406     one and three minutes, and can be used on a standard desktop computer (specification

407     of 2.3 GHz processor and 16GB memory). Typewriter, as it requires de novo assembly,

408     took up to three hours and used cloud computing or high-capacity servers.

409

410     Genefinder was written by MD. It used a mapping approach (similar to SRST2,

411     https://github.com/katholt/srst2) to detect the presence or absence of predefined

412     genes or variations in predefined genes using Bowtie. Thresholds were defined at 90%

413     overall, but amended where required in order to distinguish between both variants

414     where genes were represented with multiple reference sequences and the level of

415     diversity expected for each gene sought. Genefinder also checked for premature stop

416     codons and compared the average depth of read coverage to identify any potential

417     sequence contamination.

418

419     Mykrobe was written by PB and ZI (9). A threshold frequency was generated for each

420     gene (K minimum percentage) based on the empirical level of diversity observed in the

421     training set described by Bradley (K=0.3 for *blaZ*, K=0.6 for *fusB*, *fusC*, K=0.8 otherwise).

422     The maximum likelihood from 3 models (gene absent, gene present in minor proportion,

423     gene present) was chosen. The models took into account expected proportion of kmers

424     based on depth of coverage and empirical level of diversity (described in (9)). Mutations

425     were genotyped by choosing the maximum likelihood model from 3 Poisson models

426     comparing the depth of coverage across 63 base pair reference and alternate alleles

427     while demanding 100% coverage across the allele, also described in (9).

428

429     Typewriter was developed by TG (described in (10)). It considered BLAST results over a

430     query reference (blastn for sequence identity, tblastn for mutations). It used a "relative

431    coverage" to determine presence/absence of a gene, a metric that gives equal weight to

432    coverage and sequence identity. Typewriter reported this value for each query gene of

433    interest and cutoffs were adjusted to optimize specificity/sensitivity for different genes.

434    In this study, a relative cutoff of 90% for resistance and toxin genes was used except

435    *blaZ* for which a cutoff of 80% was used. For variant reporting, mutations were reported

436    above a given threshold of relative coverage (e.g. 90%) however, this could be changed

437    or set to 0% to report all identified differences from the query sequence. Stop codons

438    were predicted, as were novel mutations.

439

440    84 genes were included in the analysis; 46 acquired resistance genes, five sets of

441    chromosomal variants within resistance-associated genes, five cassette chromosome

442    recombinases (*ccr*) and 28 virulence genes (Tables 2, 3). Acquired resistance genes were

443    classified as present (p,P) or absent (a, A), setting 3 missing Genefinder predictions

444    ("ND" or "X") to absent. Chromosomal resistance variants were those listed in

445    Supplementary Table 4; 23 other mutations were reported in the relevant genes but

446    were not compared, as they are not considered resistance-determinants

447    (Supplementary Table 4). For all methods, genotype predictions of susceptibility

448    phenotype were based on the presence of any relevant resistance-determinant as shown

449    in Tables 2 and 3 (as described in (10) with minor modifications and updates from (9)).

450    Intermediate phenotype results were excluded from analysis (80 cases; 0.5%).

451

452

453

454

455

456    REFERENCES

457    1.    Lowy FD. 1998. Staphylococcus aureus infections. N Engl J Med 339:520-32.

458    2.    Public Health England. 2016. Annual epidemiological commentary: mandatory

459          MRSA, MSSA and E. coli bacteraemia and C. difficile infection data 2015/16

460          (https://www.gov.uk/government/uploads/system/uploads/attachment_data/f

461          ile/535635/AEC_final.pdf).

462    3.    Cosgrove SE, Qi Y, Kaye KS, Harbarth S, Karchmer AW, Carmeli Y. 2005. The

463          impact of methicillin resistance in Staphylococcus aureus bacteremia on patient

464          outcomes: mortality, length of stay, and hospital charges. Infect Control Hosp

465          Epidemiol 26:166-74.

466    4.    Banerjee R, Teng CB, Cunningham SA, Ihde SM, Steckelberg JM, Moriarty JP, Shah

467          ND, Mandrekar JN, Patel R. 2015. Randomized Trial of Rapid Multiplex

468          Polymerase Chain Reaction-Based Blood Culture Identification and Susceptibility

469          Testing. Clin Infect Dis 61:1071-80.

470    5.    Strauss C, Endimiani A, Perreten V. 2015. A novel universal DNA labeling and

471          amplification system for rapid microarray-based detection of 117 antibiotic

472          resistance genes in Gram-positive bacteria. J Microbiol Methods 108:25-30.

473    6.    Berthet N, Dickinson P, Filliol I, Reinhardt AK, Batejat C, Vallaeys T, Kong KA,

474          Davies C, Lee W, Zhang S, Turpaz Y, Heym B, Coralie G, Dacheux L, Burguiere AM,

475          Bourhy H, Old IG, Manuguerra JC, Cole ST, Kennedy GC. 2008. Massively parallel

476          pathogen identification using high-density microarrays. Microb Biotechnol 1:79-

477          86.

478    7.    Price JR, Didelot X, Crook DW, Llewelyn MJ, Paul J. 2013. Whole genome

479          sequencing in the prevention and control of Staphylococcus aureus infection. J

480          Hosp Infect 83:14-21.

481    8.     Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O,

482           Aarestrup FM, Larsen MV. 2012. Identification of acquired antimicrobial

483           resistance genes. J Antimicrob Chemother 67:2640-4.

484    9.     Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, Earle S, Pankhurst LJ,

485           Anson L, de Cesare M, Piazza P, Votintseva AA, Golubchik T, Wilson DJ, Wyllie DH,

486           Diel R, Niemann S, Feuerriegel S, Kohl TA, Ismail N, Omar SV, Smith EG, Buck D,

487           McVean G, Walker AS, Peto TE, Crook DW, Iqbal Z. 2015. Rapid antibiotic-

488           resistance predictions from genome sequence data for Staphylococcus aureus

489           and Mycobacterium tuberculosis. Nat Commun 6:10063.

490    10.    Gordon NC, Price JR, Cole K, Everitt R, Morgan M, Finney J, Kearns AM, Pichon B,

491           Young B, Wilson DJ, Llewelyn MJ, Paul J, Peto TE, Crook DW, Walker AS, Golubchik

492           T. 2014. Prediction of Staphylococcus aureus antimicrobial resistance by whole-

493           genome sequencing. J Clin Microbiol 52:1182-91.

494    11.    Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH, Del Ojo Elias C, Johnson JR,

495           Walker AS, Peto TE, Crook DW. 2013. Predicting antimicrobial susceptibilities for

496           Escherichia coli and Klebsiella pneumoniae isolates using whole genomic

497           sequence data. J Antimicrob Chemother 68:2234-44.

498    12.    Zankari E, Hasman H, Kaas RS, Seyfarth AM, Agerso Y, Lund O, Larsen MV,

499           Aarestrup FM. 2013. Genotyping using whole-genome sequencing is a realistic

500           alternative to surveillance based on phenotypic antimicrobial susceptibility

501           testing. J Antimicrob Chemother 68:771-7.

502    13.    McDermott PF, Tyson GH, Kabera C, Chen Y, Li C, Folster JP, Ayers SL, Lam C, Tate

503           HP, Zhao S. 2016. Whole-Genome Sequencing for Detecting Antimicrobial

504           Resistance in Nontyphoidal Salmonella. Antimicrob Agents Chemother 60:5515-

505           20.

506    14.    Dhup V, Kearns AM, Pichon B, Foster HA. 2015. First report of identification of

507          livestock-associated MRSA ST9 in retail meat in England. Epidemiol Infect

508          143:2989-92.

509    15.    Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, Grundman

510          H, Hasman H, Holden M, Hopkins KL, Iredell J, Kahlmeter G, Koser CU, MacGowan

511          A, Mevius D, Mulvey M, Naas T, Peto T, Rolain JM, Samuelsen O, Woodford N.

512          2016. The Role of Whole Genome Sequencing (WGS) in Antimicrobial

513          Susceptibility Testing of Bacteria: Report from the EUCAST Subcommittee. Clin

514          Microbiol Infect doi:10.1016/j.cmi.2016.11.012.

515    16.    Inouye M, Dashnow H, Raven LA, Schultz MB, Pope BJ, Tomita T, Zobel J, Holt KE.

516          2014. SRST2: Rapid genomic surveillance for public health and hospital

517          microbiology labs. Genome Med 6:90.

518    17.    Leopold SR, Goering RV, Witten A, Harmsen D, Mellmann A. 2014. Bacterial

519          whole-genome sequencing revisited: portable, scalable, and standardized

520          analysis for typing and detection of virulence and antibiotic resistance genes. J

521          Clin Microbiol 52:2365-70.

522    18.    Koser CU, Holden MT, Ellington MJ, Cartwright EJ, Brown NM, Ogilvy-Stuart AL,

523          Hsu LY, Chewapreecha C, Croucher NJ, Harris SR, Sanders M, Enright MC, Dougan

524          G, Bentley SD, Parkhill J, Fraser LJ, Betley JR, Schulz-Trieglaff OB, Smith GP,

525          Peacock SJ. 2012. Rapid whole-genome sequencing for investigation of a neonatal

526          MRSA outbreak. N Engl J Med 366:2267-75.

527    19.    Aanensen DM, Feil EJ, Holden MT, Dordel J, Yeats CA, Fedosejev A, Goater R,

528          Castillo-Ramirez S, Corander J, Colijn C, Chlebowicz MA, Schouls L, Heck M,

529          Pluister G, Ruimy R, Kahlmeter G, Ahman J, Matuschek E, Friedrich AW, Parkhill J,

530          Bentley SD, Spratt BG, Grundmann H, European SRLWG. 2016. Whole-Genome

531          Sequencing for Routine Pathogen Surveillance in Public Health: a Population

532          Snapshot of Invasive Staphylococcus aureus in Europe. MBio 7.

533   20.   Strauss L, Ruffing U, Abdulla S, Alabi A, Akulenko R, Garrine M, Germann A,

534          Grobusch MP, Helms V, Herrmann M, Kazimoto T, Kern W, Mandomando I, Peters

535          G, Schaumburg F, von Muller L, Mellmann A. 2016. Detecting Staphylococcus

536          aureus Virulence and Resistance Genes: a Comparison of Whole-Genome

537          Sequencing and DNA Microarray Technology. J Clin Microbiol 54:1008-16.

538   21.   Munson SH, Tremaine MT, Betley MJ, Welch RA. 1998. Identification and

539          characterization of staphylococcal enterotoxin types G and I from Staphylococcus

540          aureus. Infect Immun 66:3337-48.

541   22.   Letertre C, Perelle S, Dilasser F, Fach P. 2003. Identification of a new putative

542          enterotoxin SEU encoded by the egc cluster of Staphylococcus aureus. J Appl

543          Microbiol 95:38-43.

544   23.   Jarraud S, Peyrat MA, Lim A, Tristan A, Bes M, Mougel C, Etienne J, Vandenesch F,

545          Bonneville M, Lina G. 2001. egc, a highly prevalent operon of enterotoxin gene,

546          forms a putative nursery of superantigens in Staphylococcus aureus. J Immunol

547          166:669-77.

548   24.   Bae T, Baba T, Hiramatsu K, Schneewind O. 2006. Prophages of Staphylococcus

549          aureus Newman and their contribution to virulence. Mol Microbiol 62:1035-47.

550   25.   Fitzgerald JR, Sturdevant DE, Mackie SM, Gill SR, Musser JM. 2001. Evolutionary

551          genomics of Staphylococcus aureus: insights into the origin of methicillin-

552          resistant strains and the toxic shock syndrome epidemic. Proc Natl Acad Sci U S A

553          98:8821-6.

554   26.   Omoe K, Hu DL, Takahashi-Omoe H, Nakane A, Shinagawa K. 2005.

555          Comprehensive analysis of classical and newly described staphylococcal

556  superantigenic toxin genes in Staphylococcus aureus isolates. FEMS Microbiol

557  Lett 246:191-8.

558  27.  Poovelikunnel T, Gethin G, Humphreys H. 2015. Mupirocin resistance: clinical

559  implications and potential alternatives for the eradication of MRSA. J Antimicrob

560  Chemother 70:2681-92.

561  28.  Paulsen IT, Brown MH, Littlejohn TG, Mitchell BA, Skurray RA. 1996. Multidrug

562  resistance proteins QacA and QacB from Staphylococcus aureus: membrane

563  topology and identification of residues involved in substrate specificity. Proc Natl

564  Acad Sci U S A 93:3630-5.

565  29.  Dale GE, Broger C, Hartman PG, Langen H, Page MG, Then RL, Stuber D. 1995.

566  Characterization of the gene for the chromosomal dihydrofolate reductase

567  (DHFR) of Staphylococcus epidermidis ATCC 14990: the origin of the

568  trimethoprim-resistant S1 DHFR from Staphylococcus aureus? J Bacteriol

569  177:2965-70.

570  30.  Byrne ME, Gillespie MT, Skurray RA. 1991. 4',4'' adenyltransferase activity on

571  conjugative plasmids isolated from Staphylococcus aureus is encoded on an

572  integrated copy of pUB110. Plasmid 25:70-5.

573  31.  International Working Group on the Classification of Staphylococcal Cassette

574  Chromosome E. 2009. Classification of staphylococcal cassette chromosome mec

575  (SCCmec): guidelines for reporting novel SCCmec elements. Antimicrob Agents

576  Chemother 53:4961-7.

577  32.  Martineau F, Picard FJ, Lansac N, Menard C, Roy PH, Ouellette M, Bergeron MG.

578  2000. Correlation between the resistance genotype determined by multiplex PCR

579  assays and the antibiotic susceptibility patterns of Staphylococcus aureus and

580  Staphylococcus epidermidis. Antimicrob Agents Chemother 44:231-8.

581 33. Holden MT, Hsu LY, Kurt K, Weinert LA, Mather AE, Harris SR, Strommenger B,

582      Layer F, Witte W, de Lencastre H, Skov R, Westh H, Zemlickova H, Coombs G,

583      Kearns AM, Hill RL, Edgeworth J, Gould I, Gant V, Cooke J, Edwards GF, McAdam

584      PR, Templeton KE, McCann A, Zhou Z, Castillo-Ramirez S, Feil EJ, Hudson LO,

585      Enright MC, Balloux F, Aanensen DM, Spratt BG, Fitzgerald JR, Parkhill J, Achtman

586      M, Bentley SD, Nubel U. 2013. A genomic portrait of the emergence, evolution,

587      and global spread of a methicillin-resistant Staphylococcus aureus pandemic.

588      Genome Res 23:653-64.

589 34. Kaase M, Lenga S, Friedrich S, Szabados F, Sakinc T, Kleine B, Gatermann SG.

590      2008. Comparison of phenotypic methods for penicillinase detection in

591      Staphylococcus aureus. Clin Microbiol Infect 14:614-6.

592 35. El Feghaly RE, Stamm JE, Fritz SA, Burnham CA. 2012. Presence of the bla(Z) beta-

593      lactamase gene in isolates of Staphylococcus aureus that appear penicillin

594      susceptible by conventional phenotypic methods. Diagn Microbiol Infect Dis

595      74:388-93.

596 36. Testing TECoAS. 2017. Breakpoints for interpretation of MICs and zone

597      diameters. Version 7.0.

598 37. Norstrom T, Lannergard J, Hughes D. 2007. Genetic and phenotypic identification

599      of fusidic acid-resistant mutants with the small-colony-variant phenotype in

600      Staphylococcus aureus. Antimicrob Agents Chemother 51:4438-46.

601 38. Schmitz FJ, von Eiff C, Gondolf M, Fluit AC, Verhoef J, Peters G, Hadding U, Heinz

602      HP, Jones ME. 1999. Staphylococcus aureus small colony variants: rate of

603      selection and MIC values compared to wild-type strains, using ciprofloxacin,

604      ofloxacin, levofloxacin, sparfloxacin and moxifloxacin. Clin Microbiol Infect 5:376-

605      378.

606 39. Piddock LJ, Jin YF, Webber MA, Everett MJ. 2002. Novel ciprofloxacin-resistant,

607 nalidixic acid-susceptible mutant of Staphylococcus aureus. Antimicrob Agents

608 Chemother 46:2276-8.

609 40. Mellmann A, Andersen PS, Bletz S, Friedrich AW, Kohl TA, Lilje B, Niemann S,

610 Prior K, Rossen JW, Harmsen D. 2017. High Interlaboratory Reproducibility and

611 Accuracy of Next-Generation Sequencing-Based Bacterial Genotyping in a Ring-

612 Trial. J Clin Microbiol doi:10.1128/JCM.02242-16.

613 41. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,

614 Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,

615 Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and

616 its applications to single-cell sequencing. J Comput Biol 19:455-77.

617 42. Votintseva AA, Bradley P, Pankhurst L, Del Ojo Elias C, Loose M, Nilgiriwala K,

618 Chatterjee A, Smith EG, Sanderson N, Walker TM, Morgan MR, Wyllie DH, Walker

619 AS, Peto TEA, Crook DW, Iqbal Z. 2017. Same-Day Diagnostic and Surveillance

620 Data for Tuberculosis via Whole-Genome Sequencing of Direct Respiratory

621 Samples. J Clin Microbiol 55:1285-1298.

622 43. Kozyreva VK, Truong CL, Greninger AL, Crandall J, Mukhopadhyay R, Chaturvedi

623 V. 2017. Validation and Implementation of Clinical Laboratory Improvements

624 Act-Compliant Whole-Genome Sequencing in the Public Health Microbiology

625 Laboratory. J Clin Microbiol 55:2502-2520.

626 44. White LO. 2000. UK NEQAS in antibiotic assays. J Clin Pathol 53:829-34.

627 45. Andrews JM. 2001. Determination of minimum inhibitory concentrations. J

628 Antimicrob Chemother 48 Suppl 1:5-16.

629 46. Eyre DW, Golubchik T, Gordon NC, Bowden R, Piazza P, Batty EM, Ip CL, Wilson

630 DJ, Didelot X, O'Connor L, Lay R, Buck D, Kearns AM, Shaw A, Paul J, Wilcox MH,

631        Donnelly PJ, Peto TE, Walker AS, Crook DW. 2012. A pilot study of rapid benchtop

632        sequencing of Staphylococcus aureus and Clostridium difficile for outbreak

633        detection and surveillance. BMJ Open 2.

634  47.     Lahuerta-Marin A, Guelbenzu-Gonzalo M, Pichon B, Allen A, Doumith M, Lavery

635        JF, Watson C, Teale CJ, Kearns AM. 2016. First report of lukM-positive livestock-

636        associated methicillin-resistant Staphylococcus aureus CC30 from fattening pigs

637        in Northern Ireland. Vet Microbiol 182:131-4.

638  48.     Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for

639        Illumina sequence data. Bioinformatics 30:2114-20.

640  49.     Lunter G, Goodson M. 2011. Stampy: a statistical algorithm for sensitive and fast

641        mapping of Illumina sequence reads. Genome Res 21:936-9.

642  50.     Zerbino DR, Birney E. 2008. Velvet: algorithms for de novo short read assembly

643        using de Bruijn graphs. Genome Res 18:821-9.

644

645

**Figure legends**

647    Figure 1: Determinant-by-determinant disagreements between methods

648    Each panel shows percentage difference in proportion of detected presence of each

649    determinant between the first method and the second.

650

651    Figure2: Antimicrobial susceptibility genotypic predictions compared to phenotype

652

653    Figure 3: Sensitivity and specificity of genotypic predictions of antimicrobial

654    susceptibility

655

656

668    **Table 1 Overview of Genefinder, Mykrobe and Typewriter methods and**

669    **requirements**

|  | Genefinder | MyKrobe (9) | Typewriter (10) |
|---|---|---|---|
| Method | Maps raw reads to list of target alleles using Bowtie | Looks for list of target alleles in de Bruijn assembly graph | Blasts list of target alleles against de novo assemblies* |
| Input | Fastq file | Fastq file | Genome assembly output (Velvet) |
| Required homology to declare gene presence/absence | >90% to target allele | Based on Kmer recovery: K is minimum percentage expected to be recovered for a gene; K = 0.3 for blaZ, K=0.6 for Fus B, C, K= 0.8 otherwise ** | >90% relative coverage (homologyXlength) (80% for *blaZ*) |
| Required homology to declare SNP | >90% to target: can be modified | 100% of 63 kmers required to call a variant present | >90% to target: can be modified |
| Prediction of stop codons in genes present | Yes | No: there is no assembly | Yes |
| Reads can be mapped to | Multiple targets | Single target | Single target |

|                          | Genefinder                                                                                                                                          | MyKrobe (9)                                                | Typewriter (10)                                                                                       |
| ------------------------ | --------------------------------------------------------------------------------------------------------------------------------------------------- | ---------------------------------------------------------- | ---------------------------------------------------------------------------------------------------- |
| Speed / processor        | 1 to 3 minutes on laptop with 2.3 GHz processor and 16GB memory†                                                                                    | 2 minutes on laptop with 2.3 GHz processor and 16GB memory | 3 hours for assemblies on cloud computational system, then few minutes for BLAST                     |
| Sequence quality control | Threshold adjusted if gene has multiple reference sequence or variable level of diversity, can detect potential contamination by comparing average depth of coverage | Can identify mixtures of difference species and same species | Thresholds for n50 and parallel reference-based mapping: nothing reported if below these thresholds |

670

671    * using blastn for sequence identity and tblast for mutations.

672

673    † Genefinder speed is relative to the number of genes present in the database

674

675

**Table 2: Predicted antibiotic susceptibility phenotype from WGS by Genefinder,**

677 **Mykrobe, Typewriter (n=1379)**

| Antibiotic | Susceptibility prediction for Genefinder, MyKrobe, Typewriter | | | | | | Discordant across methods (n, %) |
|---|---|---|---|---|---|---|---|
| | **RRR** | **SSS** | **RRS** | **RSR** | **RSS** | **SRS** | |
| Ciprofloxacin | 304 | 1072 | 0 | 2 | 0 | 1 | 3 (0.2%) |
| Clindamycin | 338 | 1024 | 7 | 0 | 0 | 10 | 17 (1.2%) |
| Erythromycin | 354 | 1011 | 6 | 0 | 0 | 8 | 14 (1.2%) |
| Fusidic acid | 151 | 1221 | 4 | 0 | 0 | 3 | 7 (0.5%) |
| Gentamicin | 76 | 1300 | 1 | 0 | 0 | 2 | 3 (0.2%) |
| Methicillin | 393 | 984 | 2 | 0 | 0 | 0 | 2 (0.1%) |
| Mupirocin | 15 | 1362 | 0 | 0 | 2 | 0 | 2 (0.1%) |
| Penicillin | 1,161 | 211 | 3 | 0 | 0 | 4 | 7 (0.5%) |
| Rifampicin | 23 | 1,354 | 0 | 1 | 0 | 1 | 2 (0.1%) |
| Tetracycline | 121 | 1,249 | 4 | 0 | 0 | 5 | 9 (0.7%) |
| Trimethoprim | 175 | 1,199 | 3 | 1 | 0 | 1 | 5 (0.4%) |
| Vancomycin | 0 | 1,379 | 0 | 0 | 0 | 0 | 0 (0.0%) |
| **Total (% of 16548)** | 3111 (18.8%) | 13,366 (80.8%) | 30 (0.2%) | 4 (0.02%) | 2 (0.01%) | 35 (0.2%) | 71 (0.4%) |

678

679

680

681　**Table 3: Predicted genotype and phenotype**

682　**(a) Antimicrobial susceptibility**

| | Antimicrobial susceptibility prediction from Genefinder, Mykrobe, Typewriter | | | | | | |
|---|---|---|---|---|---|---|---|
| **Laboratory phenotype** | RRR | SSS | RRS | RSR | RSS | SRS | Total |
| R | 2720 | 89 | 9 | 3 | 0 | 4 | 2825 |
| S | 97 | 11504 | 13 | 1 | 2 | 22 | 11639 |
| **Total** | 2817 | 11593 | 22 | 4 | 2 | 26 | 14464 |

683

684　**(b) Virulence genes, *ccr* genes and *mecA/mecC***

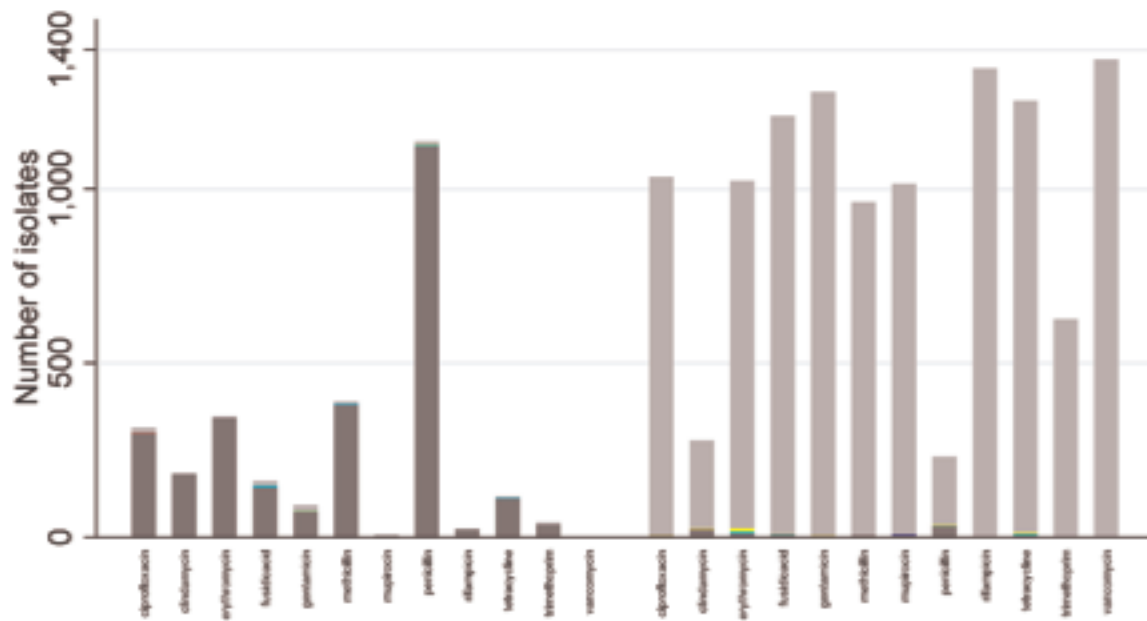| | Prediction from Genefinder, Mykrobe, Typewriter | | | | |
|---|---|---|---|---|---|
| **PCR** | AAA | PPP | APA | PPA | Total |
| A | 3362 | 82 | 10 | 17 | 3475 |
| P | 14 | 618 | 2 | 10 | 643 |
| **Total** | 3376 | 700 | 12 | 27 | 4115 |

685

686　Note: not all isolates were phenotyped for all antimicrobials, and therefore total with

687　phenotypes (14464) is less than the total with genotypic predictions (16548) in Table 2.

688　Only PHE isolates had PCR results for some virulence genes. Dark grey shading shows

689　complete concordance, and light grey majority concordance between predictions.

690　R=resistant, S=susceptible, A=absent, P=present

691

Mykrobe vs. Typewriter agreement

Genefinder vs. Typewriter agreement

Genefinder vs. Mykrobe agreement

Gene

- Genes with aquired resistance
- Genes with virulence
- Genes with chromosonal mutation
- Genes with CCR

Figure showing "Number of isolates" on the y-axis (0 to 1,400) with bar charts grouped under "Gold standard R" and "Gold standard S".

Results given as Genefinder Mykrobe Typewriter

| Color | Code |
|---|---|
| rrr | |
| rrs | |
| rsr | |
| rss | |
| srs | |
| sss | |