**The clinical utility of two high-throughput 16S rRNA gene sequencing workflows for taxonomic assignment of unidentifiable bacterial pathogens in MALDI-TOF MS**

Hiu-Yin LAO[a], Timothy Ting-Leung NG[a], Ryan Yik-Lam WONG[b], Celia Sze-Ting WONG[b], Chloe Toi-Mei CHAN[a], Denise Sze-Hang WONG[a], Lam-Kwong LEE[a] Stephanie Hoi-Ching JIM[a], Jake Siu-Lun LEUNG[a], Hazel Wing-Hei LO[a], Ivan Tak-Fai WONG[a], Miranda Chong-Yee YAU[b], Jimmy Yiu-Wing LAM[b], Alan Ka-Lun WU[b], Gilman Kit-Hang SIU[a#]

[a] *Department of Health Technology and Informatics, The Hong Kong Polytechnic University, Hong Kong Special Administrative Region, China*

[b] *Department of Clinical Pathology, Pamela Youde Nethersole Eastern Hospital, Hong Kong Special Administrative Region, China*

# Correspondence author: Dr. Gilman KH SIU, Department of Health Technology and Informatics, Hong Kong Polytechnic University, Hong Kong (gilman.siu@polyu.edu.hk)

**ABSTRACT**

Bacterial pathogens that cannot be identified using matrix-assisted laser desorption/ionization time-of-flight mass spectrometry (MALDI-TOF MS) are occasionally encountered in clinical laboratories. The *16S rRNA* gene is often used for sequence-based analysis to identify these bacterial species. Nevertheless, traditional Sanger sequencing is laborious, time-consuming and low-throughput. Here, we compared two commercially available *16S* rRNA gene sequencing tests, which are based on Illumina and Nanopore sequencing technologies, respectively, in their ability to identify the species of 172 clinical isolates that failed to be identified by MALDI-TOF MS. Sequencing data were analyzed by respective built-in analysis programs (MiSeq Reporter Software and Epi2me) and BLAST+ (v2.11.0). Their agreement with Sanger sequencing on species-level identification was determined. Discrepancies were resolved by whole-genome sequencing. The diagnostic accuracy of each workflow was determined using the composite sequencing result as the reference standard. Despite the high base-calling accuracy of Illumina sequencing, we demonstrated that the Nanopore workflow had a comparatively higher taxonomic resolution at the species level. Using built-in analysis algorithms, the concordance of Sanger 16S with the Illumina and Nanopore workflows was 33.14% and 87.79%, respectively. The agreement was 65.70% and 83.14%, respectively, when BLAST+ was used for analysis. Compared with the reference standard, the diagnostic accuracy of optimized Nanopore 16S was 96.36%, which was identical to Sanger 16S and was better than Illumina 16S (71.52%). The turnaround time of the Illumina workflow and the Nanopore workflow was 78h and 8.25h, respectively. The per-sample cost of the Illumina and Nanopore workflows was US$28.5 and US$17.7, respectively.

42

## INTROUDUCTION

44 Traditionally, clinical microbiology laboratories have relied on phenotypic methods to identify

45 bacterial pathogens. However, conventional biochemical tests are labor-intensive and time-

46 consuming, and the results can be ambiguous when two species share similar biochemical

47 profiles (1, 2). Nowadays, matrix-assisted laser desorption/ionization time-of-flight mass

48 spectrometry (MALDI-TOF MS) is widely used for bacterial identification in clinical

49 laboratories (3). MALDI-TOF MS allows rapid identification of microorganisms by comparing

50 the mass spectrum of a sample with the reference spectra in the database (4). Although MALDI-

51 TOF MS is a rapid, simple and high-throughput technology for bacterial identification, some

52 species cannot be well differentiated due to high similarity in the mass spectra of closely related

53 species or lack of reference spectra (5).

54 A study from Lau *et al*. reported that MALDI-TOF MS failed to determine the species of over

55 70% of phenotypically "difficult-to-identify" bacteria in clinical laboratories(6). In general,

56 anaerobes, particularly *Actinomyces* spp., *Peptostreptococcus* spp., *Prevotella* spp. and

57 *Fusobacterium* spp. (7-9), have a higher failure rate compared with aerobes in bacterial

58 identification using MALDI-TOF MS (7, 10). Additionally, some Gram-positive aerobes, such

59 as *Nocardia* spp. and *Streptomyces* spp., are poorly identified by MALDI-TOF MS (7, 11).

60 Regarding Gram-negative aerobes, studies show that MALDI-TOF MS cannot effectively

61 identify *Acinetobacter* spp., *Chryseobacterium* spp. and *Moraxella* spp. at the species level (11,

62 12). In such cases, *16S* sequencing of cultured isolates is commonly used for species-level

63 identification.

64  Sanger sequencing offers a high base-calling accuracy, but it is laborious and time-consuming

65  with limited throughput (13). High-throughput sequencing (HTS) technologies have been

66  proposed as alternatives to generate *16S* sequences for rapid identification of bacteria that are of

67  clinical interest. Next-generation sequencing (NGS), such as can be achieved using Illumina

68  platforms, can generate vast quantities of accurate sequencing reads. However, the read length is

69  limited and insufficient to cover the entire *16S* rRNA gene. According to the official workflow

70  for *16S* rRNA sequencing developed by Illumina Ltd., bacteria are identified based on variable

71  regions (V3 and V4) of *16S*. Nevertheless, these regions are not equally discriminative between

72  and across different species, genera and families (14).

73  The MinION device by Oxford Nanopore Technologies (ONT) enables generation of reads

74  exceeding 30 kb. The official *16S* rRNA sequencing assay allows the entire *16S* rRNA gene to

75  be sequenced with real-time data analysis. Recent studies have demonstrated its potential for

76  rapid bacterial identification; however, the high read-error rate (8%–15%) of this platform might

77  hinder the accuracy of species-level identification for diagnostic purposes (15).

78  Considering the respective limitations of Illumina and Nanopore technologies, a comprehensive

79  investigation of the clinical utility of these *16S* rRNA sequencing approaches for bacterial

80  identification is required. This study aimed to evaluate the performance of two commercial HTS

81  workflows for *16S* rRNA sequencing, namely the 16S Metagenomic Sequencing Library

82  Preparation workflow (Nextera XT Index kit v2) from Illumina and the 16S Barcoding Kit 1-24

83  (SQK-16S024) from ONT, coupled with the respective built-in analysis programs and in-house

84  BLAST+ (v2.11.0) analysis. These workflows were used to identify bacterial isolates that could

85  not be differentiated by MALDI-TOF MS. In light of the complexities of evaluating diagnostic

86  accuracy in the absence of a perfect gold standard, we considered a composite *16S* rRNA

87    sequencing result inferred by Sanger and the two HTS platforms as a reference standard. In case

88    of disagreement in taxa inferred by the three sequencing platforms, whole-genome sequencing

89    (WGS) was conducted to confirm the bacterial identities. In addition, the cost and time-to-result

90    of the sequencing workflows were also compared.

91

92    MATERIALS AND METHODS

93    **Sample collection and preparation**

94    A total of 172 clinical isolates from 117 species were collected from the clinical microbiology

95    laboratory of Pamela Youde Nethersole Eastern Hospital. Clinical isolates were included if they

96    failed to be classified at the species level (score < 2.00) by the IVD MALDI Biotyper (Bruker

97    Daltonics, Bremen, Germany). Failure to identify bacterial species occurred due to (i) lack of a

98    reference spectrum in the database (81 samples); (ii) inclusion of certain species in the

99    "dangerous database," named Security Library 1.0, rather than the regular database (two

100   samples); or (iii) poor-quality samples (89 samples) (Table S1). The IVD MALDI Biotyper used

101   in this study was microflex® (Bruker Daltonics), and the database version was BD-6763.

102   Total nucleic acid was extracted from clinical isolates using the AMPLICOR® Respiratory

103   Specimen Preparation Kit (Roche, Basel, Switzerland) and purified with 1.8X AMPure XP beads

104   (Beckman Coulter, California, USA). Purified DNA was diluted to targeted concentrations in

105   subsequent sequencing workflows. The required DNA input for the Illumina and Nanopore

106   workflows was 12.5 ng and 10 ng, respectively.

107

108   **Sanger *16S* rRNA sequencing (Sanger *16S*)**

109 The full-length *16S* rRNA gene was amplified using primers for 16s_008F (5´-

110 AGAGTTTGATCMTGGC-3´) and 16s_1507R (5´-TACCTTGTTACGACTT-3´) (16). The

111 reaction mixture was prepared by mixing 36.7 µl of nuclease-free water, 5 µl of 10× polymerase

112 chain reaction (PCR) buffer, 1 µl of 10-mM deoxynucleoside triphosphate mix (NEB, Ipswich,

113 Massachusetts, USA), 1 µl of each 25-µM primer, 0.3 µl of HotStarTaq Plus DNA Polymerase

114 (Qiagen, Hilden, Germany) and 5 µl of DNA template. The PCR conditions were 96°C for 8

115 min, 37 cycles at 94°C for 1 min, 37°C for 2 min and 72°C for 2 min 30 s, followed by 72°C for

116 10 min, and a hold step at 4°C. PCR products were purified using ExoSAP-IT reagent (Thermo

117 Fisher Scientific, Waltham, MA, USA) and then passed to the subsequent cycle sequencing using

118 eight sequencing primers (17-19) (Table S2). The reaction mixture consisted of 13 µl of

119 nuclease-free water, 1 µl of BigDye® Terminator v3.1 Ready Reaction Mix (Thermo Fisher

120 Scientific), 3.5 µl of 5× sequencing buffer, 1 µl of 3.2-µM primer and 1.5 µl of purified PCR

121 product. The PCR conditions were 96°C for 1 min, 25 cycles at 96°C for 10 sec, 37°C for 30 sec

122 and 60°C for 4 min, followed by a hold step at 4°C. The sequencing products were purified using

123 75% isopropanol and resuspended in 12 µl of Hi-Di™ Formamide (Thermo Fisher Scientific).

124 After loading on the Applied Biosystems® 3130 Genetic Analyzer (Thermo Fisher Scientific),

125 the resulting raw trace files were analyzed using the Staden Package (v2.0.0b11). The consensus

126 sequence of each sample was classified by submitting a Basic Local Alignment Search Tool

127 (BLAST) query against the *16S* ribosomal RNA sequence database.

128

129 **Illumina sequencing (NGS *16S*)**

130 *Library preparation*. Libraries were constructed according to the 16S Metagenomic Sequencing

131 Library Preparation workflow from Illumina. Briefly, the *16S* V3 and V4 regions of samples

132    were amplified in the first stage of PCR using the primers suggested in the workflow, which

133    were 16S Amplicon PCR Forward Primer (5´-

134    <u>TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG</u>CCTACGGGNGGCWGCAG-3´) and

135    16S Amplicon PCR Reverse Primer (5´-

136    <u>GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG</u>GACTACHVGGGTATCTAATCC-

137    3´). The underlined bases in the primer sequences are the overhang adapter sequences for

138    attachment of the indexed adapters in the second stage of PCR. The size of the amplicon was

139    approximately 460 bp. After a post-PCR clean-up, a unique indexed sequencing adapter was

140    added to each sample using the Nextera XT Index kit v2 (Illumina, San Diego, California, USA).

141    Then, a second post-PCR clean-up was performed, followed by a qualification check of the

142    purified libraries.

143    ***Quantification and sequencing***. The size of each library was measured using the 2100

144    Bioanalyzer system (Agilent, Santa Clara, CA, USA) and the High Sensitivity DNA kit

145    (Agilent). The quantity of the libraries was measured by real-time PCR using the LightCycler[®]

146    480 Instrument II (Roche) and QIAseq[™] Library Quant Assay Kit (Qiagen). Then, the libraries

147    were diluted to 4 nM and pooled into one tube. After denaturation with 0.2-N NaOH, the pooled

148    library was diluted to 9 pM and spiked with 15% of 9-pM PhiX prepared from PhiX Control Kit

149    v3 (Illumina). The pooled library was then loaded on the MiSeq sequencer (Illumina) for

150    sequencing using MiSeq Reagent Kits v3 (Illumina). The sequencing time was 56 h.

151    ***On-instrument data analysis***. Sequencing data were analyzed using MiSeq Reporter software

152    (v2.6.2.3) (MSR) in the MiSeq system. After selecting the metagenomics workflow, sequencing

153    reads were mapped against reference sequences in the Greengenes database (v13.5, May 2013)

154    (http://greengenes.lbl.gov/) for classification. The classification of reads at seven taxonomic

155    levels from kingdom to species was analyzed in this workflow.

156    *Data analysis using NGS_BLAST+*. The paired-end reads of each sample were merged using

157    the "make.contigs" command in Mothur (v1.44.3) (20). The reads were filtered using the

158    "screen.seqs" command. Sequences smaller than 400 bp, larger than 500 bp, or with any

159    ambiguous bases were removed. The resulting fasta files were analyzed by BLAST+ (v2.11.0)

160    using          an          in-house          Python          script

161    (https://github.com/siupenyau/Pocket_16S/tree/7d3fa9d73a6a35afb47e40e7850cef72b4b91a22).

162    In brief, the reads were aligned to the reference sequences in the *16S* ribosomal RNA database

163    (https://ftp.ncbi.nlm.nih.gov/blast/db/) downloaded from the National Center for Biotechnology

164    Information (NCBI). The percentage identity and percentage query coverage were set at 90%.

165

166    **Nanopore sequencing (Nanopore *16S*)**

167    *Library preparation and sequencing*. Library preparation was performed using the 16S

168    Barcoding Kit 1-24 (SQK-16S024) from ONT according to the manufacturer's protocol.

169    Libraries were quantified using the Qubit 2.0 Fluorometer (Thermo Fisher Scientific) with the

170    Qubit™ 1X dsDNA HS Assay Kit (Thermo Fisher Scientific). Then, 24 barcoded libraries were

171    pooled into one tube in equal concentrations. After ligation with the rapid adapter, sequencing

172    was performed using the flow cell FLO-MIN106 R9.4.1 with the MinION sequencer on the

173    MinKNOW platform for approximately 4 h.

174    *On-instrument real-time data analysis*. During sequencing, the passed fastq files, which had a

175    quality score of >7, were uploaded on the cloud-based data analysis platform Epi2me for

176 analysis. Sequencing reads were aligned to reference sequences in the NCBI 16S bacterial

177 database using the FASTQ 16S workflow (v2020. 04. 06). Regarding the workflow parameters,

178 the minimum QSCORE was set at 7, while the minimum percentage coverage and minimum

179 percentage identity were set at 90%.

180 ***Data analysis using NanoBLAST+***. In addition to Epi2me, sequencing data were analyzed using

181 BLAST+ (v2.11.0), similar to the analysis of NGS data. As each sample generated multiple fastq

182 files in a sequencing run, the fastq files of each sample were first merged into a single fastq file

183 and then converted to a fasta file before being aligned to reference sequences in the database.

184 ***Data analysis using NanoCLUST***. Samples with disagreement between EPI2ME and

185 NanoBLAST+ were further analyzed using another pipeline, NanoCLUST

186 (https://github.com/genomicsITER/NanoCLUST) (21). Unlike Epi2me and NanoBLAST+,

187 NanoCLUST does not classify individual reads in a sample. Instead, NanoCLUST forms clusters

188 of similar reads and classifies the consensus sequence of each cluster.

189 **Whole genome sequencing (WGS)**

190 Samples with complete discordant taxa, as inferred by Sanger *16S*, NGS *16S* and Nanopore *16S*

191 tests, were subjected to WGS to confirm the definite identities using the ONT platform. Library

192 preparation was performed using the transpose-based rapid barcoding kit (SQK-RBK110.96)

193 according to the manufacturer's protocol. After pooling and adapter ligation, the library was

194 loaded on the flow cell FLO-MIN106 R9.4.1 and sequenced using the GridION device for 48 h

195 in high-accuracy base calling mode. The passed fastq files were uploaded to Epi2me and

196 analyzed using the WIMP workflow (v2021.03.05).

197 **De novo assembly for WGS datasets**

198   Sequencing reads of each sample were assembled using Shasta (v0.7.0)

199   (https://github.com/chanzuckerberg/shasta). Sequencing reads were re-aligned to the assembled

200   consensus sequences using minimap2 (v2.17-r941) and samtools (v1.10). Consensus sequences

201   were first polished using MarginPolish (v1.3.dev-5492204) (https://github.com/UCSC-nanopore-

202   cgl/MarginPolish) and then further polished using homopolish (v0.2.1)

203   (https://github.com/ythuang0522/homopolish) (22). To avoid bioinformatic bias in de novo

204   assembly, each sample was also subjected to a second analysis pipeline. In brief, the sequencing

205   reads were assembled using miniasm (v0.3-r179)

206   (https://github.com/lh3/miniasm/releases/tag/v0.3). All-vs-all read self-mapping was performed

207   using minimap2. Raw consensus sequences were then generated using miniasm. After re-

208   alignment of the raw reads to consensus sequences using minimap2, the consensus sequences

209   were polished twice using racon (v1.4.3) (https://github.com/isovic/racon).

210   The longest polished consensus sequences of each sample were classified using BLAST+

211   (v2.11.0) with the Prokaryotic RefSeq Genomes database downloaded from the NCBI. The top

212   classified species with both query coverage and percentage identity were reported. The average

213   nucleotide identity (ANI) between the query and best-matched reference genomes was calculated

214   using an ANI calculator (https://www.ezbiocloud.net/tools/ani) (23). ANI >94% indicated that

215   the samples belong to the same species as the best-matched genomes.

216   **Data and statistical analysis**

217   The top classified taxa obtained from NGS and Nanopore datasets were compared with those

218   inferred by Sanger *16S* using built-in programs and BLAST+ for analysis. Species-level

219   concordance between the HTS and Sanger workflows was calculated. For samples that did not

220   match at the species level, concordance at the genus or family level was determined.

221    To assess diagnostic accuracy, a composite *16S* rRNA sequencing result of the three sequencing

222    platforms was considered as the reference standard. Identical species obtained by at least two

223    sequencing platforms were considered as reference taxa. For samples with complete discordant

224    species inferred by the three sequencing platforms, WGS was conducted to confirm the reference

225    taxa.

226

227

228    RESULTS

229    **Statistics of sequencing reads generated from the NGS and Nanopore workflows**

230    Based on the default analysis of MSR, the NGS platform generated an average of 113,381 reads

231    per sample. After merging the paired-end reads and filtering out unwanted reads with undesired

232    read lengths and ambiguous bases, an average of 68,652 filtered reads per sample was retained

233    for NGS_BLAST+ analysis.

234    The Nanopore MinKNOW platform generated an average of 51,769 reads (QSCORE ≥ 7) per

235    sample, but an average of 51,419 reads (QSCORE ≥ 7) per sample was analyzed in the FASTQ

236    16S workflow in Epi2me. The slight difference in the number of average reads per sample was

237    due to using different algorithms in the demultiplexing step between Epi2me and Guppy

238    (MinKNOW). An average of 51,769 reads per sample was analyzed using NanoBLAST+.

239    The total number of reads and the number of classified reads of each sample on both sequencing

240    platforms are shown in Table S3.

241

242    **Taxonomic resolution of sequencing reads**

243    The percentage distribution of classified reads via both sequencing platforms is shown in Figure

244    1. On average, only 45.74% of the total reads of a sample were successfully classified at the

245    species level by MSR with reference to the Greengenes database. After merging paired-end reads

246    and quality filtering, 94.02% of filtered reads were classified at the species level by

247    NGS_BLAST+ with reference to the NCBI *16S* rRNA database.

248    In the Nanopore workflow, both Epi2me and NanoBLAST+ use the NCBI *16S* rRNA database

249    for classification of long-read sequencing data. An average of 76.03% of total reads were

250    classified at the species level in Epi2me, compared with 53.56% in NanoBLAST+.

251

## Concordance in bacterial speciation by Sanger, Illumina and Nanopore *16S* rRNA sequencing

254    The top-ranked species obtained from the NGS *16S* and Nanopore *16S* workflows, coupled with

255    the respective analysis pipelines, are listed in Table S3 The percentage of samples that matched

256    with Sanger *16S* at each of the species, genus and family levels is illustrated in Figure 2. The

257    concordance in species-level identification among the sequencing platforms is shown in Figure

258    3. Overall, in terms of concordance with the Sanger *16S* result, Nanopore *16S* was better than

259    NGS *16S* (154/172 [89.53%] vs. 113/172 [65.70%], respectively), regardless of analysis

260    pipeline.

261    For the NGS *16S* workflow, MSR and NGS_BLAST+ demonstrated a concordance of 33.14%

262    (57/172) and 65.70% (113/172), respectively, with Sanger *16S* in species-level identification. A

263    total of 9.30% of samples (16/172) were unmatched, even at the family level, in MSR, whereas

264    all samples matched at the family level or below in NGS_BLAST+. Of note, concordance

265    between the results of MSR and NGS_BLAST+ was low; only 32.56% of samples (56/172)

266    showed a matched result among the classified species from these two analysis pipelines.

267    Moreover, only 28.49% of samples (49/172) showed complete agreement in the classified

268    species among the MSR, NGS_BLAST+, and Sanger datasets. Owing to poor concordance of

269    the MSR analysis with other sequencing methods, NGS_BLAST+ was considered as the optimal

270    analysis method for the Illumina datasets, and its results were regarded as the final identification

271    inferred by the NGS *16S* workflow.

272    For Nanopore *16S*, a concordance of 87.79% (151/172) and 83.14% (143/172) at the species

273    level was achieved with Epi2me and NanoBLAST+, respectively. A total of 1.16% of samples

274    (2/172) were unmatched, as reported by Epi2me and NanoBLAST+, respectively. Concordance

275    between the results of Epi2me and NanoBLAST+ was 80.23% (138/172). Additionally, 76.74%

276    of samples (132/172) showed agreement in the classified species among the Epi2me,

277    NanoBLAST+ and Sanger datasets.

278    A total of 34 samples showed disagreement in the classified species inferred by Epi2me and

279    NanoBLAST+. The respective Nanopore data were further analyzed using NanoCLUST to

280    resolve the discrepancies. NanoCLUST agreed with Epi2ME and BLAST+ in 13 (38.24%) and

281    17 (50.00%) samples, respectively. Four samples failed to reach agreement in terms of species-

282    level identification, in which three were matched in terms of genus-level identification, and one

283    was considered as having no reliable bacterial ID. Concordance between the resolved Nanopore

284    16S and Sanger 16S was 89.53% (154/172).

285

286    **WGS for bacterial isolates with discrepant species-level ID**

287    Eight samples (4.65% [8/172]) showed complete discordance in bacterial species, as inferred by

288    the three *16S* rRNA sequencing workflows. WGS was conducted to identify definite taxa.

289    Interestingly, seven of these samples failed to match with the published bacterial genomes, with

290    query coverage of <70% for the longest consensus sequences (Table 1). The ANIs to the best-

291    matched genomes were <85% (Threshold for the same species should be >94%) , suggesting that

292    these seven "difficult-to-identify" isolates were likely novel bacterial species. As the definite

293    bacterial species could not be confirmed, these samples were excluded from the subsequent

294    diagnostic evaluation.

295    The consensus sequence of one sample (R062) showed an overall query coverage of >92%, with

296    99.17% identity to *Klebsiella michiganensis* (NZ_CP060111.1). As the ANI achieved 98.71%,

297    *K. michiganensis* was therefore considered as the reference taxon for this sample.

298

299    **Diagnostic accuracy of the three *16S* rRNA sequencing workflows**

300    Considering the composite of *16S* rRNA sequencing and WGS results as reference standards, the

301    diagnostic accuracy of Sanger *16S*, NGS *16S* and Nanopore *16S* was 96.36% (159/165), 71.52%

302    (118/165) and 96.36% (159/165), respectively, for species-level identification of "difficult-to-

303    identify" bacterial pathogens (Figure 3). The mismatched samples in at least one of the

304    sequencing methods were listed in Table 2.  The diagnostic performance of each sequencing

305    workflow was summarized in Table 3.

306

307    **Comparison of sample-to-report time and running cost of the two HTS technologies**

308    The Illumina platform enables sequencing of up to 384 samples per run, whereas, owing to the

309    limited choice of sequencing barcodes, the Nanopore platform can only support a batch of 24

310    samples per run. Without considering the time for DNA extraction, it took 78 h for the Illumina

311    workflow to generate sequencing data for each run (Figure 4). With the Nanopore platform, the

312    sequencing workflow required 8.25 h. Of note, although base-calling and Epi2me analyses are

313  real-time processes, their speed is highly dependent on the strength of the computer. However,

314  Nanopore sequencing can be stopped once sufficient reads have been generated.

315  The running cost of the Nanopore workflow is relatively lower than that of the Illumina

316  workflow. The cost of the Illumina workflow per sequencing run is US $4,931 (172 samples),

317  and the cost per sample is approximately US $28.7. If the sample size is increased to 384, the

318  cost of the Illumina workflow per sequencing run is US $8,279; therefore, the cost per sample is

319  reduced to US $21.6. For the Nanopore workflow, the cost per sequencing run (24 samples) is

320  US $424, which means that the cost per sample is approximately US $17.7.

321

322    **DISCUSSION**

323    Although the majority of bacterial pathogens can be identified by MALDI-TOF MS, *16S* rRNA

324    gene sequencing is needed in clinical microbiology laboratories to confirm the identities of

325    "difficult-to-identify" clinical isolates. With reduced costs, simplified protocols and automated

326    bioinformatics pipelines, HTS has been proposed as a better alternative to Sanger sequencing for

327    sequence-based bacterial identification in clinical laboratories. This is the first study to compare

328    the performance (and evaluate the clinical utility) of two commercially available high-throughput

329    *16S* rRNA gene sequencing assays with built-in analysis software for taxonomic assignment of

330    bacterial pathogens that are unidentifiable using MALDI-TOF MS.

331    With the Illumina platform, the concordance of the classified species between MSR and Sanger

332    *16S* was exceptionally low; only 33.14% of samples matched the reference at top classified

333    species compared with 65.70% when using NGS_BLAST+. As described in previous studies, the

334    use of different bioinformatic tools and *16S* rRNA sequence databases could result in different

335    taxonomic assignments, especially at lower taxonomic levels (24, 25). The latest version of the

336    Greengenes database for MSR was updated in 2013 and does not contain certain new bacterial

337    taxa, which accounts for the poor agreement of this workflow compared with others (25).

338    Nevertheless, mismatches between NGS and Sanger sequencing were observed in 34.33% of

339    samples, even when the same aligner (i.e., BLAST+) and database (i.e., NCBI 16S bacterial

340    database) were used. One may argue that, with the constraint of low sequencing depth, the

341    Sanger *16S* result alone should not be considered as the final reference. We used a composite of

342    *16S* sequencing results generated by three platforms, and any discrepancies were resolved by

343    WGS as the reference standard to determine the diagnostic accuracy of the HTS workflows.

344    Eventually, a total of 47 samples, including 29 genera and 37 species (Table S3), remained

345 discordant between NGS *16S* and the reference standard. As indicated by Johnson *et al*.,

346 although some sub-regions (e.g., V1–V3) of *16S* s rRNA gene provide a reasonable

347 approximation of *16S* diversity, most do not capture sufficient sequence variation to discriminate

348 between closely related taxa. Also, different sub-regions show bias in the bacterial taxa that can

349 be identified (26). In this study, V3–V4 regions might perform poorly in classifying the genera

350 of discordant samples.

351 Availability of third-generation technologies means that it is becoming possible to exploit the

352 full discriminatory potential of the entire *16S* rRNA gene in a high-throughput manner. The

353 Nanopore *16S* workflow demonstrated a considerably higher percentage concordance with the

354 Sanger *16S* workflow compared with the NGS *16S* workflow, regardless of the analysis pipeline

355 used. In contrast to the built-in analysis on the Illumina platform (i.e., MSR), the performance of

356 Epi2me with Nanopore *16S* was comparable to that of nanoBLAST+ (83.14%), with 87.79% of

357 samples matching Sanger *16S* at top classified species.

358 Notably, species-level disagreement between Epi2me and nanoBLAST+ was observed in 34

359 samples (19.77%) and was subsequently resolved by NanoCLUST. Epi2me and BLAST+ rely

360 on read-by-read alignment to reference sequences in the database. As the base-calling accuracy

361 of Nanopore sequencing is relatively low, the prevalence of sequencing errors in Nanopore reads

362 could limit its ability to resolve highly similar sequences. Alternatively, NanoCLUST generates

363 clusters based on Uniform Manifold Approximation and Projection and classifies the

364 representative consensus read in each cluster using BLAST. The effect of sequencing errors in

365 individual sequences can be minimized by forming clusters, which reduces the chance of

366 misclassification. Comparing the species resolved using NanoCLUST with the reference

367    standard, there was a slight improvement in diagnostic accuracy from 89.09% (Epi2me) and

368    89.70% (nanoBLAST+) to 96.36%.

369    Six samples (3.64%) failed to match the reference at the species level in the optimized Nanopore

370    *16S* workflow. One possible reason for this discordance is the high similarity in *16S* rRNA gene

371    sequences between the inferred species and the reference taxa. Based on the now historic

372    assumption of *16S* rRNA sequencing, sequences with >95% identity represent the same genus,

373    whereas sequences with >97% identity represent closely related species (27). Many researchers

374    have reported that the taxonomic resolution of *16S* rRNA gene is lower and is unable to

375    discriminate the closely related species in certain genera, including but not limited to *Bacillus,*

376    *Burkholderia, Acinetobacter baumannii-calcoaceticus complex, Achromobacter, Actinomyces*

377    and *Staphylococcus* and the Enterobacteriaceae family (28, 29). In this study, all six taxa inferred

378    by Nanopore *16S* had >97% sequence identity with the reference standard (Table 2).

379    In this study, WGS was performed to identify the definite bacterial taxa for samples with

380    completely discordant *16S* results. To validate the transposase-based rapid sequencing protocol

381    for bacterial genome construction, two reference strains, namely *Klebsiella pneumoniae*

382    BAA3079 and *Staphylococcus aureus* BAA3114, were sequenced and analyzed in parallel with

383    the eight discordant samples. Both strains successfully yielded consensus sequences of >3Mb,

384    which covered 94% of the genomes of the respective target organisms with 99% identity. This

385    indicated that the WGS protocol was able to construct reliable consensus prokaryotic genomes

386    (Table 1). Nonetheless, the longest consensus sequences of the seven discordant samples failed

387    to obtain a query coverage >50% when mapped to the NCBI Prokaryotic RefSeq Genomes

388    database, suggesting no significant matches between these samples and published bacterial

389    genomes. The ANIs to the best-matched genomes were <94%. These "difficult-to-identify"

390 isolates were therefore considered as novel bacterial species (30). WGS confirmed that R062

391 belonged to *K. michiganensis* (ANI = 98.71%), which shared a high degree of *16S* rRNA identity

392 with the taxa assigned by Sanger *16S* (*Klebsiella grimontii*; 99.20%), NGS *16S* (*Enterobacter*

393 *cloacae*; 97.07%) and Nanopore *16S* (*Yokenella regensburgei*; 98.56%) (Table 1). This explains

394 why *16S* rRNA sequencing was not able to accurately differentiate these species.

395 Considering the time-to-result of the two sequencing platforms, the Nanopore workflow has a

396 much shorter turnaround time compared with the Illumina workflow (8.25 h and 78 h,

397 respectively). Therefore, faster results can be obtained with the Nanopore workflow. However,

398 the sample size is limited to 24 samples per batch. Comparing the cost per sample in a

399 sequencing run, Nanopore sequencing is relatively cheaper than Illumina sequencing (US $17.7

400 vs. US $28.6, respectively). Additionally, the startup cost of Nanopore sequencing is remarkably

401 lower than that of Illumina sequencing. The starter package of Nanopore sequencing costs only

402 US $1,000, whereas Illumina MiSeq costs approximately US $125,000.

403 The reusable flow cell FLO-MIN106 R9.4.1, which enables sequencing for up to 72 h, was used

404 for Nanopore 16S in this study. However, library carry over from previous run was observed in a

405 pilot study. This is problematic when the same barcode set is used in consecutive sequencing

406 run. To avoid contamination by library carry over, a new flow cell was used in each sequencing

407 run, and used flow cells were reserved for other sequencing runs using different barcodes. In this

408 context, the disposable Flongle flow cell from ONT is more suitable in a clinical setting. The

409 Flongle flow cell, which costs only US $90, can sequence for up to 16 h. Although the number of

410 active pores available in the Flongle flow cell is lower, it is more cost- and time-effective when

411 the sample size is small. Since it takes time to accumulate a batch of 24 "difficult-to-identify"

412 isolates in clinical laboratories, a small sample size per sequencing run will be beneficial,

413 especially for cases that require urgent diagnosis.

414 There are some limitations in this study that should be noted. First, the aim of this study was to

415 compare commercially available kits for *16S* rRNA gene sequencing from Illumina and

416 Nanopore. Therefore, by using the 16S Metagenomic Sequencing Library Preparation kit, only

417 the V3–V4 sub-regions of *16S* rRNA gene were sequenced in the Illumina workflow. But it is

418 possible to sequence full-length *16S* rRNA gene using Ilumina MiSeq with a laboratory

419 developed protocol(31), which may increase the taxonomic resolution of the Illumina workflow

420 at the species level. Second, except for the eight discordant samples, the reference taxa of

421 isolates were defined by *16S* rRNA sequencing without being confirmed by WGS. However,

422 some closely related species may have identical *16S* rRNA genes; thus, *16S* rRNA sequencing

423 results may not represent the definite taxa of these samples. Third, regarding the eight samples

424 that underwent WGS, the taxonomic assignment was based on the contigs of consensus

425 sequences after de novo assembly. Circular, gap-free bacterial genomes were not constructed.

426 Finally, bacterial DNA for *16S* sequencing was extracted from cultured isolates. The

427 performance of the NGS *16S* and Nanopore *16S* workflows on direct bacterial identification in

428 microbial and polymicrobial specimens was not evaluated.

429

430 **CONCLUSION**

431 In conclusion, the commercial *16S* rRNA gene sequencing workflow from ONT (SQK-16S024),

432 coupled with NanoCLUST, is the most accurate for bacterial identification in a clinical setting,

433    with higher flexibility in sample size and sequencing time, a lower running cost, and higher

434    concordance with the reference standard.

435

440    **DECLARATION OF INTEREST STATEMENT**

441    We declare no competing interests.

442

443    **REFERENCES**

444    1.    Jesumirhewe C, Ogunlowo PO, Olley M, Springer B, Allerberger F, Ruppitsch W. 2016.
445          Accuracy of conventional identification methods used for Enterobacteriaceae isolates in
446          three Nigerian hospitals. PeerJ 4:e2511.
447    2.    Harmsen D, Rothganger J, Frosch M, Albert J. 2002. RIDOM: Ribosomal Differentiation
448          of Medical Micro-organisms Database. Nucleic Acids Res 30:416-7.
449    3.    Karger A. 2016. Current developments to use linear MALDI-TOF spectra for the
450          identification and typing of bacteria and the characterization of other cells/organisms
451          related to infectious diseases. Proteomics Clin Appl 10:982-993.
452    4.    Patel R. 2015. MALDI-TOF MS for the diagnosis of infectious diseases. Clin Chem
453          61:100-11.
454    5.    Hou TY, Chiang-Ni C, Teng SH. 2019. Current status of MALDI-TOF mass
455          spectrometry in clinical microbiology. J Food Drug Anal 27:404-414.
456    6.    Lau SK, Tang BS, Teng JL, Chan TM, Curreem SO, Fan RY, Ng RH, Chan JF, Yuen
457          KY, Woo PC. 2014. Matrix-assisted laser desorption ionisation time-of-flight mass
458          spectrometry for identification of clinically significant bacteria that are difficult to
459          identify in clinical laboratories. J Clin Pathol 67:361-6.
460    7.    Ge MC, Kuo AJ, Liu KL, Wen YH, Chia JH, Chang PY, Lee MH, Wu TL, Chang SC, Lu
461          JJ. 2017. Routine identification of microorganisms by matrix-assisted laser desorption
462          ionization time-of-flight mass spectrometry: Success rate, economic analysis, and clinical
463          outcome. J Microbiol Immunol Infect 50:662-668.
464    8.    Garner O, Mochon A, Branda J, Burnham CA, Bythrow M, Ferraro M, Ginocchio C,
465          Jennemann R, Manji R, Procop GW, Richter S, Rychert J, Sercia L, Westblade L,

466    Lewinski M. 2014. Multi-centre evaluation of mass spectrometric identification of
467    anaerobic bacteria using the VITEK(R) MS system. Clin Microbiol Infect 20:335-9.
468  9.   Knoester M, van Veen SQ, Claas EC, Kuijper EJ. 2012. Routine identification of clinical
469    isolates of anaerobic bacteria: matrix-assisted laser desorption ionization-time of flight
470    mass spectrometry performs better than conventional identification methods. J Clin
471    Microbiol 50:1504.
472  10.  Luo Y, Siu GKH, Yeung ASF, Chen JHK, Ho PL, Leung KW, Tsang JLY, Cheng VCC,
473    Guo L, Yang J, Ye L, Yam WC. 2015. Performance of the VITEK MS matrix-assisted
474    laser desorption ionization-time of flight mass spectrometry system for rapid bacterial
475    identification in two diagnostic centres in China. J Med Microbiol 64:18-24.
476  11.  Bizzini A, Jaton K, Romo D, Bille J, Prod'hom G, Greub G. 2011. Matrix-assisted laser
477    desorption ionization-time of flight mass spectrometry as an alternative to 16S rRNA
478    gene sequencing for identification of difficult-to-identify bacterial strains. J Clin
479    Microbiol 49:693-6.
480  12.  Homem de Mello de Souza HAP, Dalla-Costa LM, Vicenzi FJ, Camargo de Souza D,
481    Riedi CA, Filho NAR, Pilonetto M. 2014. MALDI-TOF: a useful tool for laboratory
482    identification of uncommon glucose non-fermenting Gram-negative bacteria associated
483    with cystic fibrosis. J Med Microbiol 63:1148-1153.
484  13.  Winand R, Bogaerts B, Hoffman S, Lefevre L, Delvoye M, Braekel JV, Fu Q, Roosens
485    NH, Keersmaecker SC, Vanneste K. 2019. Targeting the 16s Rrna Gene for Bacterial
486    Identification in Complex Mixed Samples: Comparative Evaluation of Second (Illumina)
487    and Third (Oxford Nanopore Technologies) Generation Sequencing Technologies. Int J
488    Mol Sci 21.
489  14.  Chakravorty S, Helb D, Burday M, Connell N, Alland D. 2007. A detailed analysis of
490    16S ribosomal RNA gene segments for the diagnosis of pathogenic bacteria. J Microbiol
491    Methods 69:330-9.
492  15.  Ip CLC, Loose M, Tyson JR, de Cesare M, Brown BL, Jain M, Leggett RM, Eccles DA,
493    Zalunin V, Urban JM, Piazza P, Bowden RJ, Paten B, Mwaigwisya S, Batty EM,
494    Simpson JT, Snutch TP, Birney E, Buck D, Goodwin S, Jansen HJ, O'Grady J, Olsen HE,
495    Min IONA, Reference C. 2015. MinION Analysis and Reference Consortium: Phase 1
496    data release and analysis. F1000Res 4:1075.
497  16.  Muyzer G, Teske A, Wirsen CO, Jannasch HW. 1995. Phylogenetic relationships of
498    Thiomicrospira species and their identification in deep-sea hydrothermal vent samples by
499    denaturing gradient gel electrophoresis of 16S rDNA fragments. Arch Microbiol
500    164:165-72.
501  17.  Liu Z, Lozupone C, Hamady M, Bushman FD, Knight R. 2007. Short pyrosequencing
502    reads suffice for accurate microbial community analysis. Nucleic Acids Res 35:e120.
503  18.  Andersson AF, Lindberg M, Jakobsson H, Backhed F, Nyren P, Engstrand L. 2008.
504    Comparative analysis of human gut microbiota by barcoded pyrosequencing. PLoS One
505    3:e2836.
506  19.  Nossa CW, Oberdorf WE, Yang L, Aas JA, Paster BJ, Desantis TZ, Brodie EL, Malamud
507    D, Poles MA, Pei Z. 2010. Design of 16S rRNA gene primers for 454 pyrosequencing of
508    the human foregut microbiome. World J Gastroenterol 16:4135-44.
509  20.  Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,
510    Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ,
511    Weber CF. 2009. Introducing mothur: open-source, platform-independent, community-

512    supported software for describing and comparing microbial communities. Appl Environ
513    Microbiol 75:7537-41.
514  21.  Rodriguez-Perez H, Ciuffreda L, Flores C. 2020. NanoCLUST: a species-level analysis
515    of 16S rRNA nanopore sequencing data. Bioinformatics
516    doi:10.1093/bioinformatics/btaa900.
517  22.  Huang YT, Liu PY, Shih PW. 2021. Homopolish: a method for the removal of systematic
518    errors in nanopore sequencing by homologous polishing. Genome Biol 22:95.
519  23.  Yoon SH, Ha SM, Lim J, Kwon S, Chun J. 2017. A large-scale evaluation of algorithms
520    to calculate average nucleotide identity. Antonie Van Leeuwenhoek 110:1281-1286.
521  24.  Sierra MA, Li Q, Pushalkar S, Paul B, Sandoval TA, Kamer AR, Corby P, Guo Y, Ruff
522    RR, Alekseyenko AV, Li X, Saxena D. 2020. The Influences of Bioinformatics Tools and
523    Reference Databases in Analyzing the Human Oral Microbial Community. Genes (Basel)
524    11.
525  25.  Park SC, Won S. 2018. Evaluation of 16S rRNA Databases for Taxonomic Assignments
526    Using Mock Community. Genomics Inform 16:e24.
527  26.  Johnson JS, Spakowicz DJ, Hong BY, Petersen LM, Demkowicz P, Chen L, Leopold SR,
528    Hanson BM, Agresta HO, Gerstein M, Sodergren E, Weinstock GM. 2019. Evaluation of
529    16S rRNA gene sequencing for species and strain-level microbiome analysis. Nat
530    Commun 10:5029.
531  27.  Schloss PD, Handelsman J. 2005. Introducing DOTUR, a computer program for defining
532    operational taxonomic units and estimating species richness. Appl Environ Microbiol
533    71:1501-6.
534  28.  Janda JM, Abbott SL. 2007. 16S rRNA gene sequencing for bacterial identification in the
535    diagnostic laboratory: pluses, perils, and pitfalls. J Clin Microbiol 45:2761-4.
536  29.  Church DL, Cerutti L, Gurtler A, Griener T, Zelazny A, Emler S. 2020. Performance and
537    Application of 16S rRNA Gene Cycle Sequencing for Routine Identification of Bacteria
538    in the Clinical Microbiology Laboratory. Clin Microbiol Rev 33.
539  30.  Konstantinidis KT, Tiedje JM. 2005. Genomic insights that advance the species
540    definition for prokaryotes. Proc Natl Acad Sci U S A 102:2567-72.

541

542

**Table 1: Whole genome sequencing analysis for the samples with complete discordant taxonomic assiagnment by Sanger, NGS and Nanopore 16s rRNA sequencing**

| | | | | Whole genome sequencing (WGS) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Genome assembly method | | | | | |
| | | | | | Shasta | | | Miniasm | | |
| Sample ID | Species inferred by Sanger 16s | Species inferred by NGS 16s | Species inferred by Nanopore 16s | Best-matched Species by WGS (reference genome) | Query coverage (%) | Identity (%) | ANI (%)[b] | Query coverage (%) | Identity (%) | ANI (%)[b] |
| *Klebsiella pneumoniae* BAA3079[a] | *Klebsiella pneumoniae* | *Klebsiella pneumoniae* | *Klebsiella pneumoniae* | *Klebsiella pneumoniae* (NC_016845.1) | 99 | 97.00 | **98.92** | 92.13 | 99.40 | **99.14** |
| *Staphylococcus aureus* BAA3114[a] | *Staphylococcus aureus* | *Staphylococcus aureus* | *Staphylococcus aureus* | *Staphylococcus aureus* (NC_007795.1) | 94.06 | 99.95 | **99.30** | 88.39 | 99.92 | **99.23** |
| R001 | *Kocuria koreensis* | *Kocuria massiliensis* | *Kocuria spp.* | *Kocuria massiliensis* (NZ_LT835161.1) | 42.21 | 87.44 | **78.29** | 42.42 | 87.41 | **78.55** |
| R006 | *Kocuria koreensis* | *Kocuria massiliensis* | *Kocuria spp.* | *Kocuria massiliensis* (NZ_LT835161.1) | 43.04 | 79.12 | **78.49** | 42.04 | 87.49 | **78.44** |
| R062 | *Klebsiella grimontii* | *Enterobacter cloacae* | *Yokenella regensburgei* | *Klebsiella michiganensis* (NZ_CP060111.1) | 92.17 | 99.17 | **98.71** | 86.30 | 98.99 | **98.69** |
| R120 | *Brachybacterium conglomeratum* | *Brachybacterium faecium* | *Brachybacterium paraconglomeratum* | *Brachybacterium saurashtrense* (NZ_CP031356.1) | 62.15 | 85.18 | **82.30** | 62.30 | 85.12 | **82.39** |
| R121 | *Schaalia odontolytica* | *Schaalia vaccimaxillae* | *Sphingomonas paucimobilis* | *Schaalia odontolytica* (NZ_CP046315.1) | 6.07 | 78.55 | **70.34** | 6.04 | 78.24 | **70.86** |
| R131 | *Schaalia odontolytica* | *Schaalia vaccimaxillae* | *No reliable ID* | *Schaalia odontolytica* (NZ_CP046315.1) | 6.19 | 82.12 | **71.21** | 6.29 | 78.25 | **71.26** |
| R158 | *Microbacterium ginsengiterrae* | *Microbacterium assamensis* | *Microbacterium foliorum* | *Microbacterium foliorum* (NZ_CP041040.1 ) | 65.41 | 84.52 | **82.24** | 65.21 | 84.51 | **82.15** |
| R181 | *Sphingomonas yabuuchiae* | *Sphingomonas paucimobilis* | *Sphingomonas sanguinis* | *Sphingomonas hominis* (NZ_JABULH010000007.1) | 31.48 | 89.67 | **82.09** | 30.68 | 89.59 | **81.95** |

[a] *Klebsiella pneumoniae* BAA3079 and *Staphylococcus aureus* BAA3114 served as QC sample, which were sequenced and analyzed in parallel with the discordant samples for WGS and bioinformatics analysis.

[b] Average Nucleotide Identity (ANI) > 94% indicated that the samples belong to the same species as the best-matched genomes.

543

**Table 2: The samples with mismatched taxa inferred by at least one sequencing platform**

| Sample ID | Species-level ID (Reference Standard) | Sanger Sequencing (Sanger 16s) | | Illumina Sequencing (NGS 16s) | | Nanopore Sequencing (Nanopore 16s) | |
|---|---|---|---|---|---|---|---|
| | | Classified species from Sanger 16s [a] | 16s Identity against the reference (%) | Classified species from NGS 16s [a] | 16s Identity against the reference (%) | Classified species from Nanopore 16s [a] | 16s Identity against the reference (%) |
| R003 | *Pseudoglutamicibacter albus* | *Pseudoglutamicibacter cumminsii* | 99.26% | *Pseudoglutamicibacter albus* | matched | *Pseudoglutamicibacter albus* | matched |
| R013 | *Microbacterium hominis* | *Microbacterium hominis* | matched | *Microbacterium aerolatum* | 97.47% | *Microbacterium hominis* | matched |
| R017 | *Microbacterium hominis* | *Microbacterium hominis* | matched | *Microbacterium aerolatum* | 97.47% | *Microbacterium hominis* | matched |
| R021 | *Microbacterium hominis* | *Microbacterium hominis* | matched | *Microbacterium aerolatum* | 97.47% | *Microbacterium hominis* | matched |
| R024 | *Bacillus idriensis* | *Bacillus idriensis* | matched | *Bacillus idriensis* | matched | *Bacillus indicus* | 97.62% |
| R025 | *Varibaculum cambriense* | *Varibaculum cambriense* | matched | *Varibaculum anthropi* | 98.50% | *Varibaculum cambriense* | matched |
| R026 | *Varibaculum cambriense* | *Varibaculum cambriense* | matched | *Varibaculum anthropi* | 98.50% | *Varibaculum cambriense* | matched |
| R036 | *Corynebacterium lowii* | *Corynebacterium lowii* | matched | *Corynebacterium bovis* | 93.29% | *Corynebacterium lowii* | matched |
| R040 | *Weissella cibaria* | *Weissella cibaria* | matched | *Weissella confusa* | 99.26% | *Weissella cibaria* | matched |
| R043 | *Proteus vulgaris* | *Proteus vulgaris* | matched | *Proteus alimentorum* | 99.64% | *Proteus vulgaris* | matched |
| R045 | *Brucella microti* | *Brucella microti* | matched | *Brucella papionis* | 99.86% | *Brucella microti* | matched |
| R047 | *Proteus cibarius* | *Proteus cibarius* | matched | *Proteus terrae* | 99.65% | *Proteus cibarius* | matched |
| R049 | *Dermacoccus barathri* | *Dermacoccus barathri* | matched | *Dermacoccus profundi* | 99.86% | *Dermacoccus barathri* | matched |
| R052 | *Arcanobacterium wilhelmae* | *Arcanobacterium wilhelmae* | matched | *Arcanobacterium pinnipediorum* | 96.60% | *Arcanobacterium wilhelmae* | matched |
| R053 | *Dermacoccus barathri* | *Dermacoccus barathri* | matched | *Dermacoccus profundi* | 99.86% | *Dermacoccus barathri* | matched |
| R056 | *Corynebacterium simulans* | *Corynebacterium simulans* | matched | *Corynebacterium glutamicum* | 93.74% | *Corynebacterium simulans* | matched |
| R058 | *Corynebacterium mastitidis* | *Corynebacterium mastitidis* | matched | *Corynebacterium tuberculostearicum* | 94.67% | *Corynebacterium mastitidis* | matched |
| R062 | *Klebsiella michiganensis* | *Klebsiella grimontii* | 99.20% | *Enterobacter cloacae* | 97.07% | *Yokenella regensburgei* | 98.56% |
| R063 | *Corynebacterium pilbarense* | *Corynebacterium pilbarense* | matched | *Corynebacterium coyleae* | 98.04% | *Corynebacterium pilbarense* | matched |
| R069 | *Eikenella corrodens* | *Eikenella corrodens* | matched | *Eikenella halliae* | 98.69% | *Eikenella corrodens* | matched |
| R071 | *Corynebacterium xerosis* | *Corynebacterium hansenii* | 99.07% | *Corynebacterium xerosis* | matched | *Corynebacterium xerosis* | matched |
| R072 | *Mycolicibacterium fortuitum* | *Mycolicibacterium fortuitum* | matched | *Mycolicibacterium arcueilense* | 98.96% | *Mycolicibacterium fortuitum* | matched |
| R073 | *Tessaracoccus oleiagri* | *Tessaracoccus oleiagri* | matched | *Tessaracoccus flavescens* | 95.95% | *Tessaracoccus oleiagri* | matched |
| R078 | *Vagococcus teuberi* | *Vagococcus teuberi* | matched | *Vagococcus martis* | 99.22% | *Vagococcus teuberi* | matched |
| R079 | *Corynebacterium xerosis* | *Corynebacterium hansenii* | 99.07% | *Corynebacterium xerosis* | matched | *Corynebacterium xerosis* | matched |
| R083 | *Tessaracoccus oleiagri* | *Tessaracoccus oleiagri* | matched | *Tessaracoccus flavescens* | 95.95% | *Tessaracoccus oleiagri* | matched |
| R086 | *Raoultella planticola* | *Raoultella planticola* | matched | *Raoultella planticola* | matched | *Klebsiella aerogenes* | 99.06% |
| R094 | *Corynebacterium xerosis* | *Corynebacterium hansenii* | 99.07% | *Corynebacterium xerosis* | matched | *Corynebacterium xerosis* | matched |

| R096 | *Streptomyces thermodiastaticus* | *Streptomyces thermodiastaticus* | matched | *Streptomyces thermoviolaceus* | 98.86% | Streptomyces thermodiastaticus | matched |
|---|---|---|---|---|---|---|---|
| R097 | *Pseudoxanthomonas helianthi* | *Pseudoxanthomonas helianthi* | matched | *Pseudoxanthomonas spadix* | 97.04% | *Pseudoxanthomonas helianthi* | matched |
| R098 | *Brachybacterium huguangmaarense* | *Brachybacterium huguangmaarense* | matched | *Brachybacterium huguangmaarense* | matched | *Brachybacterium nesterenkovii* | 97.84% |
| R104 | *Gordonia sputi* | *Gordonia sputi* | matched | *Gordonia otitidis* | 99.07% | *Gordonia sputi* | matched |
| R105 | *Gordonia sputi* | *Gordonia sputi* | matched | *Gordonia otitidis* | 99.07% | *Gordonia sputi* | matched |
| R108 | *Staphylococcus saccharolyticus* | *Staphylococcus saccharolyticus* | matched | *Staphylococcus epidermidis* | 99.19% | *Staphylococcus saccharolyticus* | matched |
| R112 | *Citrobacter sedlakii* | *Citrobacter sedlakii* | matched | *Citrobacter youngae* | 98.32% | *Citrobacter sedlakii* | matched |
| R116 | *Tsukamurella tyrosinosolvens* | *Tsukamurella tyrosinosolvens* | matched | *Tsukamurella ocularis* | 99.86% | *Tsukamurella tyrosinosolvens* | matched |
| R123 | *Pseudoglutamicibacter albus* | *Pseudoglutamicibacter cumminsii* | 99.26% | *Pseudoglutamicibacter albus* | matched | *Pseudoglutamicibacter albus* | matched |
| R133 | *Nocardia brasiliensis* | *Nocardia brasiliensis* | matched | *Nocardia vulneris* | 99.31% | *Nocardia brasiliensis* | matched |
| R140 | *Moraxella lacunata* | *Moraxella lacunata* | matched | *Moraxella equi* | 99.38% | *Moraxella lacunata* | matched |
| R141 | *Ottowia beijingensis* | *Ottowia beijingensis* | matched | *Brachymonas denitrificans* | 93.33% | *Ottowia beijingensis* | matched |
| R149 | *Ornithinibacillus californiensis* | *Ornithinibacillus californiensis* | matched | *Ornithinibacillus scapharcae* | 98.48% | *Ornithinibacillus californiensis* | matched |
| R151 | *Dermacoccus barathri* | *Dermacoccus barathri* | matched | *Dermacoccus profundi* | 99.86% | *Dermacoccus barathri* | matched |
| R153 | *Corynebacterium mastitidis* | *Corynebacterium mastitidis* | matched | *Corynebacterium tuberculostearicum* | 94.67% | *Corynebacterium mastitidis* | matched |
| R175 | *Corynebacterium pollutisoli* | *Corynebacterium pollutisoli* | matched | *Corynebacterium humireducens* | 98.07% | *Corynebacterium pollutisoli* | matched |
| R176 | *Tsukamurella ocularis* | *Tsukamurella ocularis* | matched | *Tsukamurella ocularis* | matched | *Tsukamurella hominis* | 100.00% |
| R178 | *Acinetobacter soli* | *Acinetobacter soli* | matched | *Acinetobacter soli* | matched | *Acinetobacter lactucae* | 97.82% |
| R179 | *Corynebacterium lipophiloflavum* | *Corynebacterium lipophiloflavum* | matched | *Corynebacterium mycetoides* | 97.16% | *Corynebacterium lipophiloflavum* | matched |
| R180 | *Corynebacterium mastitidis* | *Corynebacterium mastitidis* | matched | *Corynebacterium tuberculostearicum* | 94.67% | *Corynebacterium mastitidis* | matched |
| R182 | *Fusobacterium nucleatum* | *Fusobacterium nucleatum* | matched | *Fusobacterium canifelinum* | 98.34% | *Fusobacterium nucleatum* | matched |
| R183 | *Parabacteroides faecis* | *Parabacteroides faecis* | matched | *Parabacteroides chongii* | 97.15% | *Parabacteroides faecis* | matched |
| R190 | *Bacillus xiamenensis* | *Bacillus xiamenensis* | matched | *Bacillus aerius* | 97.16% | *Bacillus xiamenensis* | matched |
| R192 | *Corynebacterium pilbarense* | *Corynebacterium pilbarense* | matched | *Corynebacterium ureicelerivorans* | 98.85% | *Corynebacterium pilbarense* | matched |
| R204 | *Prevotella scopos* | *Prevotella scopos* | matched | *Prevotella jejuni* | 97.41% | *Prevotella scopos* | matched |
| R205 | *Pasteurella multocida* | *Pasteurella multocida* | matched | *Pasteurella stomatis* | 93.74% | *Pasteurella multocida* | matched |
| R206 | *Staphylococcus cohnii* | *Staphylococcus cohnii* | matched | *Staphylococcus auricularis* | 98.16% | *Staphylococcus cohnii* | matched |
| R208 | *Achromobacter denitrificans* | *Achromobacter denitrificans* | matched | *Achromobacter xylosoxidans* | 99.15% | *Achromobacter denitrificans* | matched |
| R210 | *Bacillus licheniformis* | *Bacillus licheniformis* | matched | *Bacillus piscis* | 97.37% | *Bacillus licheniformis* | matched |

544        [a] The mismatched taxa were underlined.

**Table 3: Diagnostic accuracies of the Sanger, NGS and Nanopore 16s rRNA sequencing methods**

| Sequencing method | No. of sample analyzed | No. of samples with matched taxa | Diagnostic Accuracy (%) | 95% CI |
|---|---|---|---|---|
| **Sanger 16s** | **165** | **159** | **96.36** | **92.25 - 98.65** |
| **Optimized NGS 16s [a]** | **165** | **118** | **71.52** | **63.98 - 78.26** |
| Analyzed by MSR | 165 | 59 | 35.76 | 28.46 - 43.58 |
| Analyzed by NGS_BLAST+ | 165 | 118 | 71.52 | 63.98 - 78.26 |
| **Optimized Nanopore 16s [b]** | **165** | **159** | **96.36** | **92.25 - 98.65** |
| Analyzed by Epi2ME | 165 | 147 | 89.09 | 83.31 - 93.41 |
| Analyzed by NanoBLAST+ | 165 | 148 | 89.7 | 84.02 - 93. 88 |

[a] Owing to the poor concordance of MSR with other methods, the NGS_BLAST+ was considered as the optimal analysis method for the Illumina datasets

[b] The mismatched taxa inferred by Epi2ME and NanoBLAST+ were resolved by NanoCLUST.

545

546

547 **FIGURE LEGENDS**



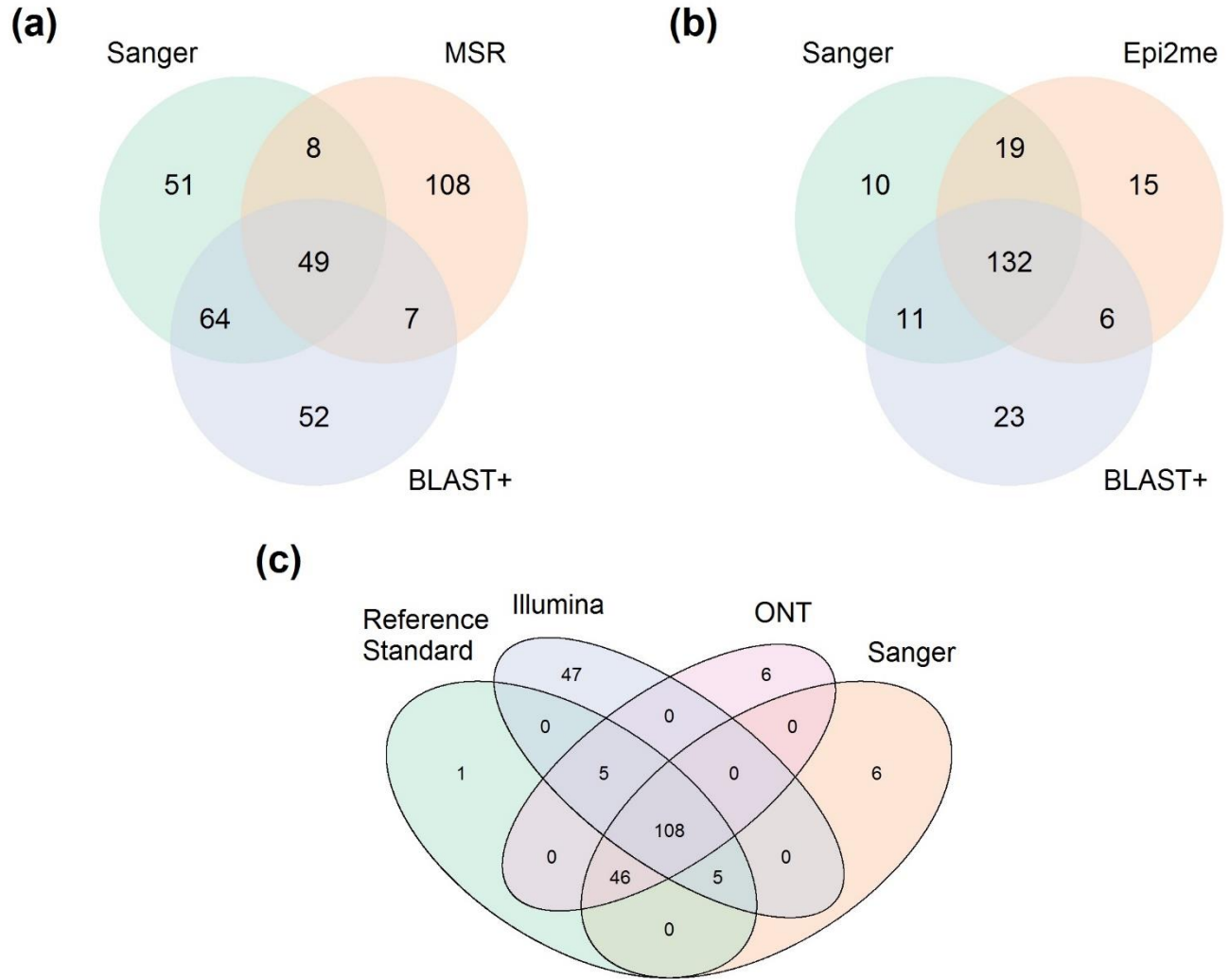549 Figure 1. The boxplots showing the distribution of percentage of classified reads of all samples in (a)

550 Illumina and (b) Nanopore sequencing.

**The concordance between bacterial taxa generated from HTS technologies and the reference Sanger sequencing**

| | Miseq reporter (Illumina) | NGS_BLAST+ (Illumina) | Epi2me (Nanopore) | NanoBLAST+ (Nanopore) |
|---|---|---|---|---|
| ■ Not match | 9.30% | 0.00% | 1.16% | 1.16% |
| ■ Matched at family level | 18.02% | 1.16% | 0.58% | 1.16% |
| ■ Matched at genus level | 39.54% | 33.14% | 10.47% | 14.54% |
| ■ Matched at species level | 33.14% | 65.70% | 87.79% | 83.14% |

Figure 2. The concordance between bacterial taxa inferred by the two HTS workflows and the Sanger sequencing.

555

Figure 3. The Venn Diagram showing the concordance of bacterial taxa inferred by different 16S rRNA sequencing platforms. (a) Concordance of top classified species between Illumina sequencing, coupled with MSR and NGS_BLAST+ analysis, and Sanger sequencing. (b) Concordance of top classified species between Nanopore sequencing, coupled with Epi2ME and nanoBLAST+, and Sanger sequencing. (c) Concordance of top classified species among Sanger 16S, NGS 16S, Nanopore 16S and reference standard.
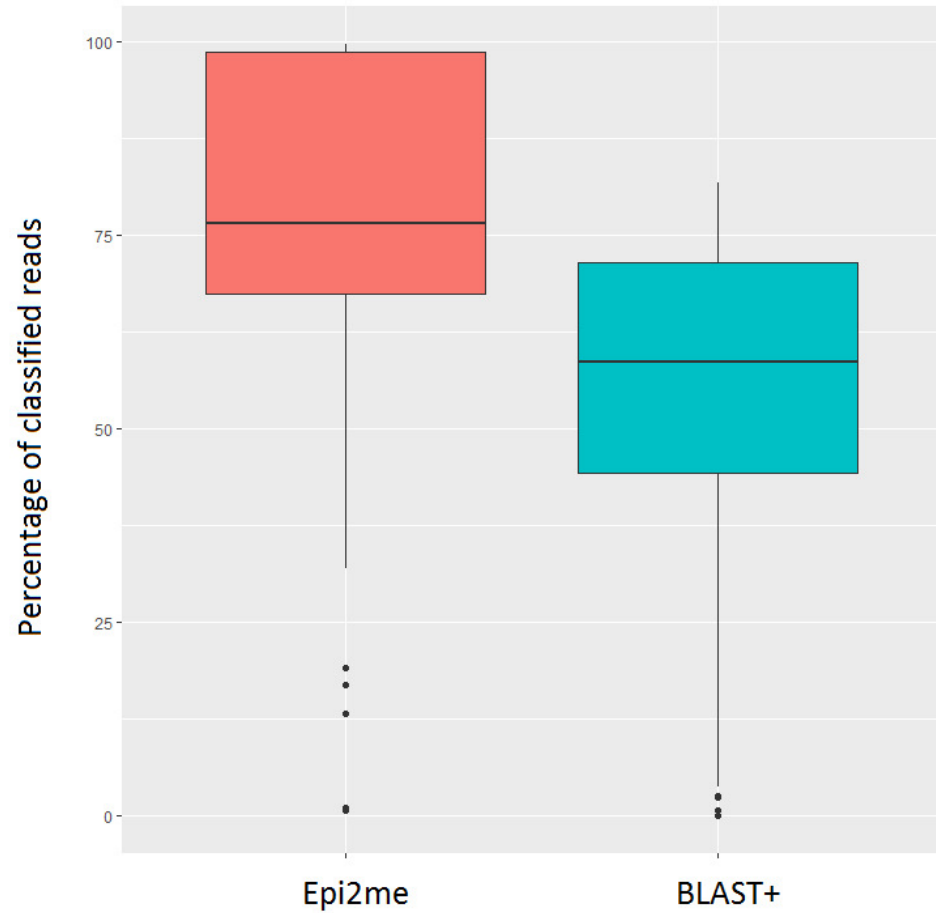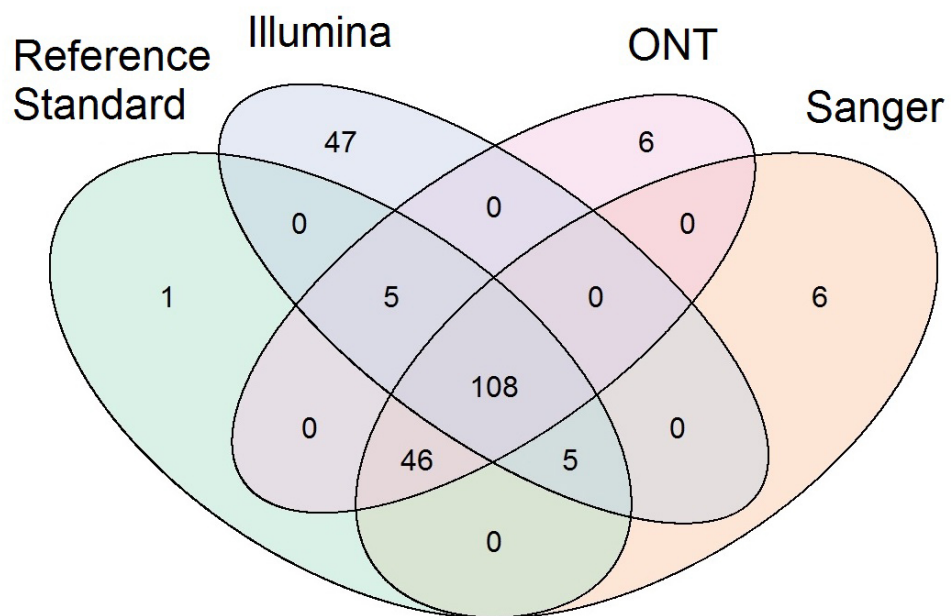
562

563

564

565

566

567

568

569

570



571

Figure 4. 16S rRNA gene sequencing workflow of the HTS technologies.

573

# The concordance between bacterial taxa generated from HTS technologies and the reference Sanger sequencing



| | Miseq reporter (Illumina) | NGS_BLAST+ (Illumina) | Epi2me (Nanopore) | NanoBLAST+ (Nanopore) |
|---|---|---|---|---|
| Not match | 9.30% | 0.00% | 1.16% | 1.16% |
| Matched at family level | 18.02% | 1.16% | 0.58% | 1.16% |
| Matched at genus level | 39.54% | 33.14% | 10.47% | 14.54% |
| Matched at species level | 33.14% | 65.70% | 87.79% | 83.14% |

**(a)**

Sanger      MSR

51    8    108

49

64    7

52

BLAST+

**(b)**

Sanger      Epi2me

10    19    15

132

11    6

23

BLAST+

**(c)**

Reference Standard    Illumina    ONT    Sanger

47    6

0    0    0

1    5    0    6

108

0    0

46    5

0