1   # Validation and Implementation of CLIA-Compliant

2   # Whole Genome Sequencing (WGS) in Public Health Laboratory

3

4

5   Varvara K. Kozyreva, Chau-Linda Truong, Alexander L. Greninger, John

6   Crandall, Rituparna Mukhopadhyay, and Vishnu Chaturvedi*

7

8   Microbial Diseases Laboratory, California Department of Public Health,

9   Richmond, CA

10

11  Bacteria, Whole Genome Sequencing, Performance Specifications,

12  Laboratory Developed Test, Quality Management, Validation, CLIA, Public

13  Health

14

15

16  *Corresponding author: Vishnu.Chaturvedi@cdph.ca.gov

## Abstract

## Background

Public health microbiology laboratories (PHL) are at the cusp of unprecedented improvements in pathogen identification, antibiotic resistance detection, and outbreak investigation by using whole genome sequencing (WGS). However, considerable challenges remain due to the lack of common standards.

## Objectives

1) Establish the performance specifications of WGS applications used in PHL to conform with CLIA (Clinical Laboratory Improvements Act) guidelines for laboratory developed tests (LDT), 2) Develop quality assurance (QA) and quality control (QC) measures, 3) Establish reporting language for end users with or without WGS expertise, 4) Create a validation set of microorganisms to be used for future validations of WGS platforms and multi-laboratory comparisons and, 5) Create modular templates for the validation of different sequencing platforms.

## Methods

MiSeq Sequencer and Illumina chemistry (Illumina, Inc.) were used to generate genomes for 34 bacterial isolates with genome sizes from 1.8 to 4.7 Mb and wide range of GC content (32.1%-66.1%). A customized CLCbio Genomics Workbench - shell script bioinformatics pipeline was used for the data analysis.

## Results

44  We developed a validation panel comprising ten *Enterobacteriaceae* isolates,

45  five gram-positive cocci, five gram-negative non-fermenting species, nine

46  *Mycobacterium tuberculosis*, and five miscellaneous bacteria; the set

47  represented typical workflow in the PHL.  The accuracy of MiSeq platform for

48  individual base calling was >99.9% with similar results shown for

49  reproducibility/repeatability of genome-wide base calling.  The accuracy of

50  phylogenetic analysis was 100%.  The specificity and sensitivity inferred

51  from MLST and genotyping tests were 100%.  A test report format was

52  developed for the end users with and without WGS knowledge.

## Conclusion

54  WGS was validated for routine use in PHL according to CLIA guidelines for

55  LDTs.  The validation panel, sequencing analytics, and raw sequences will be

56  available for future multi-laboratory comparisons of WGS in PHL.

57  Additionally, the WGS performance specifications and modular validation

58  template are likely to be adaptable for the validation of other platforms and

59  reagents kits.

60

3

## Introduction

Clinical and public health microbiology laboratories are undergoing transformative changes with the adoption of whole genome sequencing (WGS) [1, 2]. For several years, leading laboratories have published proof-of-concept studies on WGS-enabled advances in the identification of pathogens, antibiotic resistance detection, and disease outbreak investigations [3-6]. The technologies also referred to as next generation sequencing (NGS) have yielded more detailed information about the microbial features than was possible using a combination of other laboratory approaches. Further developments of WGS platforms had allowed remarkable in-depth inquiry of pathogenic genomes for the discovery of genetic variants and genome rearrangements that could have been missed using other DNA methods [3, 7, 8]. The enhanced investigations of disease outbreaks have led to new understanding of transmission routes of infectious agents [9-11]. WGS-enabled metagenomics and microbiome discoveries have revealed a new appreciation for the role of microbes in health and disease [12-15]. The innovations are continuing at such an unprecedented pace that WGS is expected to become an alternative to culture-dependent approaches in the clinical and public microbiology laboratories [16-18].

Notwithstanding its promises, several challenges remain for the adoption of WGS in microbiology laboratories [19-22]. The accelerated

83    obsolescence of the sequencing platforms presents several obstacles in

84    bridging the gap between research and routine diagnostics including

85    standardizations efforts [23]. The downstream bioinformatics pipelines are

86    also unique challenges for the microbiology laboratory both in terms of

87    infrastructure and skilled operators [24-27]. Overall, WGS 'wet bench-dry

88    bench' workflow represents an integrated process, which is not easily

89    amenable to the traditional quality metrics used by the microbiology

90    laboratories [27-29]. The capital investments and recurring costs of WGS for

91    clinical laboratories although rapidly declining still remain relatively high to

92    allow multi-laboratory comparisons for the standardization of the analytical

93    parameters. Finally, the regulatory agencies have not yet proposed WGS

94    standard guidelines for the clinical microbiology [30], and external

95    proficiency testing programs are still in development for the clinical and

96    public health microbiology laboratories [31, 32].

97

98        There are other notable recent developments towards standardization

99    and validation of next generation sequencing in clinical laboratories. The US

100   Centers for Disease Control and Prevention (CDC) sponsored the Next-

101   generation Sequencing: Standardization of Clinical Testing (Nex-StoCT)

102   workgroup to propose quality laboratory practices for the detection of DNA

103   sequence variations associated with heritable human disorders [33, 34]. The

104   workgroup developed principles and guidelines for test validation, quality

5

105    control, proficiency testing, and reference materials.  Although not focused

106    on infectious diseases, these guidelines provide a valuable roadmap for the

107    implementation of WGS in clinical microbiology and public health

108    laboratories. The College of American Pathologists' (CAP) published eighteen

109    requirements in an accreditation checklist for the next generation

110    sequencing analytic ('wet bench') and bioinformatics ('dry bench') processes

111    as part of its' molecular pathology checklist [30]. These 'foundational'

112    accreditation requirements were designed to be broadly applicable to the

113    testing of inheritable disorders, molecular oncology, and infectious diseases.

114    Along the same lines, the feasibility of *in silico* proficiency testing has been

115    demonstrated for NGS [35].  Clinical and Laboratory Standards Institute

116    (CLSI) has updated its' "Nucleic acid sequencing methods in diagnostic

117    laboratory medicine" guidelines with considerations specific to the

118    application of next generation sequencing in microbiology [36]. Thus, a

119    broad technical framework is now available to design WGS validation

120    protocols that will be most relevant for the clinical and public health

121    laboratories. Our aims for the current study were to establish performance

122    metrics for the workflow typical in the microbiological public health

123    laboratories, design modular templates for the validation of different

124    platforms and chemistries, finalize user-friendly report format, and identify a

125    set of bacterial pathogens that could be used for WGS validation and

126    performance assessments.

**Methods**

**Bacterial isolates and sequences**

A set of 34 bacterial isolates representing typical workflow in the PHL, was used for validation and quality control of WGS.  These included ten *Enterobacteriaceae* isolates, five gram-positive cocci bacterial pathogens, five gram-negative non-fermenting bacterial pathogens, nine *Mycobacterium tuberculosis* isolates and five miscellaneous bacterial pathogens (Table 1). This Whole Genome Shotgun project has been deposited at GenBank under the accession MTFS00000000-MTGZ00000000. The version described in this paper is version MTFS01000000-MTGZ01000000. Raw and assembled sequences are available for download (see Supplementary Table 1 for the accession numbers).

**Reference whole genomes**

The genome sequences of ATCC strains, isolates characterized by CDC, and other representative isolates were downloaded from NCBI database (http://www.ncbi.nlm.nih.gov/genome/) to be used as reference per the recommendations in the  CLSI guidelines [36], (Table 1).

**WGS wet bench workflow**

The whole genome sequencing was performed on Illumina MiSeq sequencer (Figure 1).  The Nextera XT library preparation procedure and 2x300 cycle MiSeq sequencing kits were used (Illumina Inc., San Diego, CA, USA).  Illumina Nextera XT indexes were used for barcoding.  Bacterial DNA

7

149    was extracted using Wizard Genomic DNA Kit (Promega, Madison, WI, USA).

150    The bacterial DNA concentrations were measured using Qubit fluorometric

151    quantitation with Qubit dsDNA BR Assay Kit (Thermo Fisher Scientific,

152    Waltham, MA, USA).  The DNA purity was estimated using NanoDrop 2000

153    UV-Vis Spectrophotometer (NanoDrop Products, Wilmington, DE, USA).  The

154    Mastercycler nexus was used for tagmentation incubation and PCR

155    (Eppendorf North America, Hauppauge, NY, USA).  The library concentration

156    was measured using Qubit HS kit.  DNA library size distribution was

157    estimated using 2100 BioAnalyzer Instrument and High Sensitivity DNA

158    analysis kit (Agilent Technologies, Santa Clara, CA, USA). Ampure beads

159    were used for size selection. Manual normalization of libraries was

160    performed.  The PhiX Control V3 sequencing control was used in every

161    sequencing run (Illumina, Inc. San Diego, CA, USA). Genomes were

162    generated with the depth coverage in the range of 15.71x-216.4x (average

163    79.72x, median 71.55x).

164    **Bioinformatics pipeline**

165     Paired-end reads were quality trimmed with the threshold of Q30, and

166    then used for mapping to the reference and *de novo* assembly on CLCbio

167    Genomic Workbench 8.0.2 (Qiagen, Aarhus, Denmark).  The BAM files

168    generated after mapping to the reference genome were taken through series

169    of software suites to generate the phylogenetic tree.  A customized shell

170    script was created to automate the subsequent steps after mapping that

171 included: 1) SNP calling in coding and non-coding genome areas using

172 SAMtools mpileup (v.1.2; [37]); 2) Converting into VCF matrix using bcftools

173 (v0.1.19; http://samtools.github.io/bcftools/); 3) Variants parsing using

174 vcftools (v.0.1.12b; [38]) to include only high-quality SNPs (hqSNPs) with

175 coverage ≥30x, minimum quality > 200; with InDels and the heterozygote

176 calls excluded. 4) Converting SNP matrix into FASTA alignment file for the

177 export back to the CLCbio GW 8.0.2 for the generation of the phylogenetic

178 tree.

179 hqSNP-based genotyping - The Maximum Likelihood phylogenetic trees were

180 generated based on high-quality single nucleotide polymorphisms (hqSNPs)

181 under the Jukes-Cantor nucleotide substitution model; with bootstrapping.

182 16S rRNA gene-based identification - Genomes were annotated with prokka

183 v1.1 tool [39] and species identification was performed by comparing 16S

184 rRNA gene sequences against the Ribosomal Database Project (RDP)

185 database [40].

186 *In silico* MLST - *In silico* multi-locus sequence typing (MLST) was performed

187 using the Center for Genomic Epidemiology (CGE) online tool [41].

188   <u>ABR genes detection</u> was performed using the CGE ResFinder online

189   resource [42]. ATCC reference strains designated for use as antibiotic

190   susceptibility controls were analyzed for the presence of antibiotic resistance

191   genes. Negative controls were chosen among strains which were described

192   by the CLSI M100-S25 document [43] as susceptible, with no known

193   antibiotic resistance genes. Positive controls were chosen among strains,

194   which according to the CLSI M100-S25 resistance determinants.

195

**Validation Plan**

196

197   Thirty-four bacterial isolates were sequenced in triplicate.  For between-run

198   reproducibility assessments, all replicates were generated starting from fresh

199   cultures except for *M. tuberculosis* where DNA samples were used.  Between

200   run replicates were processed on separate days by different operators.  For

201   within run replicates, one DNA extract was used, but independent library

202   preparations were done, with final samples being included in one sequencing

203   run.

204 **Results**

205 A number of CLIA-required quality parameters were adopted with some

206 modification for validation on WGS (Table 2).  The modular validation

207 template and a summary of performed here WGS validation for 34 bacterial

208 isolates are presented in figure 2.

209 **Accuracy of WGS**

210 The accuracy of WGS was divided into three components: platform

211 accuracy, assay accuracy, and bioinformatics pipeline accuracy.

212 Platform accuracy - Platform accuracy was assessed as the accuracy of

213 identification of individual base pairs in the bacterial genome. The accuracy

214 of the platform was established by determining the proximity of agreement

215 between base calling made by MiSeq sequencer (measured value) and NCBI

216 reference sequence (the true value). We determined MiSeq Illumina platform

217 accuracy by mapping generated reads to the corresponding reference

218 sequence and identifying Single Nucleotide Polymorphisms (SNPs).  Few

219 validation samples differed from reference genome by several SNPs.

220 However, 99% (324 out of 327) of those SNPs were reproducible among all

221 five replicates we have sequenced for each sample. Since sequencing errors

222 are random between different library preparations and it is unlikely that the

223 same erroneous SNP will occur in all 5 replicates, we can conclude that those

224 discrepancies were not caused by sequencing errors, but most likely were a

225 result of accumulation of mutations in the reference strains or previous

11

226    sequencing mistakes in the reference sequence.  In both cases, whether we

227    take into the account all SNPs detected between validation and reference

228    sequence, or only those SNPs which don't appear in all of the replicates (true

229    sequencing errors), we observed > 99.999% agreement of generated whole

230    genome sequences with the reference sequences for each tested sample.

231        Assay accuracy - Assay accuracy was determined by an agreement of

232    the assay result for the validation samples with the assay result for

233    reference sequences of the same strains. Four applications of WGS were

234    used to validate the accuracy of the assay: *in silico* Multilocus Sequence

235    Typing (MLST) assay, 16S rRNA gene species identification (ID) assay, an

236    assay for detection of antibiotic resistance (ABR) genes, and genotyping

237    assay using high-quality Single Nucleotide Polymorphisms (hqSNPs).

238        The definition of the correct result for MLST corresponds to a correct

239    identification of each of the MLST alleles in the validation sequence. For all

240    validation samples each of the sequences of the seven housekeeping genes

241    used in the typing scheme (or 6 genes- for *Aeromonas hydrophilia*) were

242    identified correctly, resulting in 100% allele identification accuracy.

243        For ABR genes detection the comparison of validation sequences was

244    performed against each entry in the ResFinder database, which at the

245    moment of validation contained sequences of 1719 antibiotic resistance

246    genes, resulting in a total of 1719 tests performed for each validation

247    sample.  In negative control samples, all 1719 tests gave negative results.

12

248   In positive controls, 1 out of 1719 tests gave a positive result, and the rest

249   must remain negative, as expected. Thus, the accuracy of the assay for ABR

250   genes detection was 100%.

251   For 16S rRNA ID assay, variations only in one gene were detected, so

252   the species ID results as a whole (e.g. "*Escherichia coli*") was considered as

253   a single test.  The identity of 16S rRNA sequence extracted from validation

254   sample showed 100% match with 16S rRNA sequence extracted from the

255   reference sequence.

256   To assess the accuracy of the genotyping test, phylogenetic trees were

257   built using reference sequences and validation sequences, and resulting

258   trees were compared. For better comparison, we used at least five strains of

259   the same species in the phylogenetic tree.  The accuracy of the genotyping

260   test was determined using two approaches: 1) Topological similarity

261   between reference tree and validation tree using Compare2Trees software,

262   and 2) Comparison of clustering pattern of validation tree and reference

263   tree.  The phylogenetic trees were generated for five bacterial isolates. All

264   five validation trees had matching clustering patterns and 100% of

265   topological similarity with corresponding reference trees (Supplementary

266   Table 2).

267   Bioinformatics pipeline accuracy - Accuracy of the bioinformatics

268   pipeline used for hqSNP genotyping was assessed by performing

13

269 phylogenetic analysis on raw WGS reads of bacterial isolates from well-

270 characterized outbreaks and comparing validation results to the previously

271 published phylogenetic results (Table 3). Two studies, presenting a

272 phylogenetic analysis of outbreaks, caused by the gram-positive pathogen in

273 one study [44] and gram-negative in another study [45] (at least six

274 isolates/study), were used for validation of the bioinformatics pipeline

275 (Figure 3).  The clustering of validation tree completely replicated clustering

276 of Study 1 [44] tree (Figure 3A-C), e.g. isolates 4 and 5 were identical and

277 clustered together according to the Study 1, and the same results were

278 shown in validation tree, with isolates 4 and 5 sharing the same node.  All

279 conclusions in regards to the genetic relatedness of the isolates that can be

280 drawn from Study 1 tree can also be made from analysis of validation tree 1.

281 The group of related isolates from Study 1 was compared with

282 epidemiologically unrelated isolates suggested by the same study (no tree

283 available from publication). The phylogenetic analysis using the PHL

284 bioinformatics pipeline showed that epidemiologically unrelated isolates did

285 not cluster with the group of outbreak isolates and appeared to be

286 genetically distant (Figure 3D). Thus, the resulting phylogenetic tree

287 produced by our bioinformatics pipeline showed complete concordance with

288 the epidemiological data.

289     From the Study 2 [45], we have selected nine isolates, which were

290 representative of 4 independent outbreaks and two isolates were

14

291    epidemiologically unrelated controls (Figure 3E-G). The clustering of

292    validation tree was identical to the clustering of Study 2 tree. For example,

293    isolates 6 and 7 were a part of the same outbreak, while isolate 8 is an

294    epidemiologically unrelated control used in the study. By epidemiological

295    data and Study 2 tree, the validation tree showed that isolates 6 and seven

296    do cluster together, but not with isolate 8. All observations about the genetic

297    relatedness of the isolates drawn from Study 2 tree could be replicated from

298    the analysis of validation tree 2.  In summary, based on analysis of

299    simulated data from both studies accuracy of the pipeline for phylogenetic

300    analysis was 100%.

301

302    **WGS repeatability and reproducibility.**

303        Repeatability (precision within run) was established by sequencing the

304    same samples multiple times under the same conditions and evaluating the

305    concordance of the assay results and performance. Reproducibility (precision

306    between runs) was assessed as the consistency of the assay results and

307    performance characteristics for the same sample sequenced on different

308    occasions.  Thirty-four validation samples each were sequenced three times

309    in the same sequencing run (for repeatability) and in 3 times in different

310    runs (for reproducibility).  Between run replicates were processed on

311    different days, altering two operators, as recommended CLSI MM11A

15

312   document [46]. For within run replicates, one DNA extract was used, but

313   independent library preparations were done, with final samples being

314   included in one sequencing run.  Therefore, for each sample, the number of

315   intra-assay replicates and inter-assay replicates were three each, and the

316   total numbers of repeated results were five.  All quality parameters [depth of

317   coverage, uniformity of coverage, and accuracy of base calling (Q score)]

318   remained relatively constant within and between runs.

319       Two methods of evaluating precision were used: evaluation of absolute

320   inter- and intra-assay precision per replicate and evaluation of precision

321   relative to the genome size. One out of 3 within-run replicates of isolate C50

322   *Pseudomonas aeruginosa* ATCC 27853 had a 1 SNP difference from other

323   within-run replicates (see Supplementary Table 3). All validation samples

324   except C50 yielded identical whole genome sequences for all three within-

325   run replicates.  The inter-assay precision was 99.02% as per replicate.

326   Three validation samples had one of the between-run replicates each

327   differing from other between-run replicates. Sample C47 *Staphylococcus*

328   *epidermidis* ATCC 12228 had one between-run replicate with 2 SNPs

329   difference from other replicates. Samples C49 *Streptococcus pneumoniae*

330   ATCC 6305 and C55 *Escherichia coli* ATCC 25922 each had one of the

331   between-run replicates differing from other replicated sequences by 1 SNP.

332   Intra-assay precision per replicate was 97.05%. If precision per base pair is

16

333    estimated (in relation to the covered genome size), both inter- and intra-

334    assay precision were > 99.9999%.

335        We also estimated reproducibility and repeatability for MLST and 16S

336    rRNA ID assays. For MLST total number of alleles analyzed for either within-

337    or between-run replicates was 441. Each single allele in all validation

338    samples was identified consistently among within- and between-run

339    replicates. Within- and between-run replicates had repeatable/reproducible

340    sequences of 16S rRNA gene and resulted in repeatable/reproducible species

341    identification. Within and between run precisions of allele detection and

342    species identification for corresponding assays were 100%.

343

344    **WGS Sensitivity and Specificity**

345        Analytical sensitivity and specificity of WGS were estimated for

346    genotyping and MLST.

347        Genotyping sensitivity and specificity - to estimate analytical sensitivity

348    and specificity of WGS-based genotyping, the hqSNPs phylogenetic trees

349    generated from the validation sequences were compared to the trees

350    generated from the reference sequences for the same strain.  All generated

351    validation trees repeated clustering and had 100% of topological similarity

352    with corresponding reference trees, indicating absence false negative or

353    false-positive results in the genotyping test.  Both analytical sensitivity and

354    analytical specificity of the hqSNP-based genotyping assay were 100%.

355    MLST sensitivity and specificity - As described above, using organism-

356    specific MLST databases sequence type of validation sequences and their

357    reference sequences was determined. For MLST number of the true positive

358    results corresponds to the number of alleles correctly identified in the

359    validation samples. For the true negative results, we performed a

360    comparison of validation sequences against MLST databases for unmatched

361    species, e.g. search of alleles for C1 *Escherichia coli* validation sample

362    against MLST database for *Salmonella enterica*.  In the latter case, the MLST

363    assay is not supposed to be able to identify any alleles.   All alleles in

364    positive validation samples were identified correctly. None of the alleles in

365    negative controls were identified. Both analytical sensitivity and analytical

366    specificity of *in silico* MLST test were 100%.

367

368    **WGS reportable range**

369    The following information about the sequenced genome was collected

370    for the reportable range: genome-wide hq SNPs, housekeeping genes used

371    in MLST schemes, 16S rRNA gene, and antibiotic resistance genes included

372    into ResFinder database.

373    Reporting language was developed to assist interpretation of the results

374    by an end user with or without specific WGS knowledge- the template and

375    examples are provided in the Supplementary Document 1.

376

**Quality assurance and quality control of WGS**

377

378    The quality assurance (QA) and quality control (QC) measures were

379    developed as the results of valuation to ensure high quality and consistency

380    of further routine testing using MiSeq Illumina platform. QC must be

381    performed during both pre-analytical (DNA isolation, library preparation),

382    analytical (quality metrics of sequencing run) and post-analytical (data

383    analysis) steps of the WGS. On the stage of data analysis, QC includes three

384    steps: raw read QC, mapping quality QC (or/and *de novo* assembly QC),

385    variant calling QC. PHL should use the WGS validation to establish the

386    thresholds of quality parameters, which can be used in following routine

387    testing to filter out poor quality samples and data and this way minimize the

388    chance of false results. We suggest spiked-in positive and negative controls

389    for routine testing as well as more comprehensive monthly positive and

390    negative controls. Since traditional CLIA rules require the positive and

391    negative control to pass through all the pre-analytical steps, including DNA

392    isolation, laboratory may choose to follow this guidance and perform DNA

393    isolation and sequencing of positive and negative control in each run, or

394    alternatively, implement Individualized Quality Control Plan (IQCP) [as per

395    42CFR493.1250] and use more economical spiked-in control instead. Type

396    and complexity of positive and negative controls should be determined by

397    each laboratory individually based on specifics of their workflow (most

398    probable source of contamination), type of microorganisms and assays which

399    are most commonly used. Regular and monthly QC practices are

400    summarized in Supplementary Figure 1. The complete QA&QC manual

401    established for WGS applications used in microbiological PHL can be found in

402    Supplementary Document 1.

403    **Validation Summary**

404      WGS assay was shown to have >99.9% accuracy, >99.9%

405    reproducibility/repeatability, and 100% specificity and sensitivity, which

406    meets CLIA requirements for laboratory-developed tests (LDTs).

407

**Discussion**

408

409     This study established the workflow and reference materials for the

410 validation of WGS for routine use in PHL according to CLIA guidelines for

411 LDTs. The validation panel, sequencing analytics, and raw sequences

412 generated during this study could serve as a resource for the future multi-

413 laboratory comparisons of WGS. Additionally, the WGS performance

414 specifications and modular validation template developed in the study could

415 be easily adapted for the validation of other platforms and reagents kits.

416 These results strengthen the concept of unified laboratory standards for

417 WGS enunciated by some professional organizations, including the Global

418 Microbial Identifier (GMI) initiative [30, 31, 33, 47]. A few other groups have

419 also highlighted the challenges and solutions for the implementation of WGS

420 in clinical and public health microbiology laboratories [21, 48].

421     Using a combination of reference strains and corresponding publicly

422 available genomes, we devised a framework of 'best practices' for the quality

423 management of the integrated 'wet lab' and 'dry lab' WGS workflow

424 ('pipeline'). The importance of reference materials for validation and QC of

425 wet- and dry-lab WGS processes has been noted earlier [28, 31, 33] . Unlike

426 in human genomics [49], there is no well-established source of reference

427 materials for WGS validation in microbiological PHL. The main challenge of

428 creating customized validation set is the lack of reference materials, in other

429 words, strains that can be easily acquired by the PHLs and which have high-

430    quality well-characterized reference genomes available. While using

431    complete genomic sequences of ATCC strains from NCBI is an option, it is far

432    from being perfect. The genome sequences available from public databases

433    are generated by using different methods, chemistries, platforms, which

434    may yield different error rates, therefore deposited sequences are not

435    guaranteed to be free of such errors. With the perpetual development of new

436    sequencing technologies and improvements in the quality of sequences, it is

437    not unlikely that the genomes sequenced with old methods may appear less

438    accurate than the validation sequences generated by the laboratory during

439    validation. In addition to this, there is a possibility of mutations

440    accumulation in the control strains, e.g. ATCC cultures, which are

441    propagated by the different laboratories. In this sense, there is no gold

442    standard available for use as a reference material for bacterial WGS

443    validation. Nevertheless, NCBI, ENA, and similar public genome depositories

444    remain to be the best resource for the genomic sequences of control strains,

445    which can be used for validation. In future, it would be optimal to have a

446    network/agency/bank which could distribute panels of thoroughly sequenced

447    isolates, with curated and updated genomic sequences available online for

448    WGS validation. In the absence of such resource, we developed a validation

449    set of microorganisms, which can be used for future validations of WGS

450    platforms. Bacterial genomes vary differently in size, GC content, abundance

451    of repetitive regions, and other properties, which affect the WGS results. We

452  created a validation set which reflects the diversity of the microorganisms

453  with various genome sizes and GC-content, which are routinely sequenced

454  by the PHL. Different species of gram-positive and gram-negative

455  microorganisms and *M. tuberculosis* were included to account for the

456  differences in DNA isolation procedures as well.

457      Samples were validated based on four core elements also reflected in

458  the assay report: 16S rRNA-based species identity, *in silico* MLST, hqSNP

459  phylogenetic analysis, and the presence of AR determinants. Overall, we

460  achieved high accuracy, precision, sensitivity and specificity for all test

461  analytes ranging from 99-100%, which well exceeds 90% threshold for

462  these performance parameters for LDT as per CLIA. These findings are in

463  agreement with recent reports of 93%-100% accuracy in WGS identification,

464  subtyping, and antimicrobial resistance genes detection in a number of

465  pathogens [50-53].

466      The successful CLIA integration of the WGS would also obligate a

467  laboratory to implement a continuous performance measurement plan via an

468  internal or external proficiency testing (PT) program. Such PT programs are

469  under active development with the Global Microbial Identifier (GMI) network,

470  the Genetic Testing Reference Materials Coordination Program (Get-RM), the

471  Genome in a Bottle (GIAB) Consortium, and the CDC PulseNet NextGen

472  being the most prominent [31, 49]. More generic standards have been

473  proposed by the College of American Pathologists' (CAP) molecular

23

474  pathology checklist (MOL)[30]. The proposed quality standards include both

475  live cultures as well as 'sequence only' formats for a comprehensive

476  assessment of the WGS pipeline. Our validation set of isolates is amenable

477  to both internal and external quality assurance testing. In preliminary

478  internal PT, we were able to successfully assess the entire workflow and

479  personnel performance (details not shown).

480       Microbial WGS remains a dynamic technology, and therefore, any

481  validated pipeline is unlikely to remain static. For this reason,

482  implementation of modular validation template becomes crucial for the

483  seamless and timely introduction of changes to the 'pipeline,' e.g. we had to

484  carry-out several amendments to the protocol since its implementation in

485  the laboratory. These included a new processing algorithm for highly-

486  contagious pathogens and some adjustments to the data analysis algorithm.

487  The changes were accomplished via minor modifications in the 'pipeline' with

488  corroborative testing using developed by us modular validation template. We

489  also performed a two-sequencer comparison to allow for processing of

490  increased volume of samples (see the protocol for the correlation study in

491  Supplementary Document 1).

492       The WGS report format continues to pose challenges. Reporting

493  language was designed to be able to convey the WGS-based assay results to

494  the end user with or without the extensive knowledge of WGS to avoid

495  erroneous interpretation of the results by the final user and provide

24

496 actionable data. Disclaimers are particularly important to guide the potential

497 use of the data in clinical settings, e.g. a disclaimer that detection of

498 antibiotic resistance genes by WGS do not guarantee resistance of the strain

499 *in vivo* and that phenotypic susceptibility test is required to confirm

500 antimicrobial resistance.

501 The study possesses certain limitations. Firstly, only a limited number

502 of WGS-based assays were included into the validation study based on the

503 most common PHL applications. Other types of WGS assays/analytics would

504 have to be validated in a similar manner to determine the performance

505 specifications, which are required to generate accurate and reproducible

506 results, e.g. a threshold for the base calling accuracy of the platform, or a

507 depth of coverage of specific genes. Secondly, not all validation set samples

508 had available NCBI database entries to provide comparison sets. Thirdly, the

509 absence of any eukaryotic pathogens in the current validation is another

510 shortcoming and therefore, additional validation studies would be needed to

511 implement a pipeline for the pathogenic fungi and parasites.

512 As the clinical and public microbiology community implements high-

513 quality WGS, it would be opportune to consider the available models for the

514 delivery of these services [54]. Since their inception, most WGS activities

515 have taken place in the reference facilities with rather large supporting

516 infrastructure. Although inevitable in the early stages, the centralization of

517 services presents several challenges on the turnaround time and access to

25

518 the specific expertise on the local population structure of a given pathogen,

519 which are crucial for the management of infectious diseases at the local and

520 regional levels. WGS services could now be delivered locally, more easily

521 with the affordable sequencers, standardized reagents, and well-defined

522 quality metrics. The local delivery model would also be more responsive to

523 the needs of the target client and enhance the adoption of WGS across the

524 healthcare systems. Another alternative is a hybrid model with

525 complimentary central and local services to balance the need for speed with

526 the advanced expertise and resources [54]. Two prominent examples of the

527 hybrid models in the United States are the Food and Drug Administration

528 (FDA) GenomeTrakr network for the tracking of food-borne pathogens, and

529 the CDC Advanced Molecular Detection (AMD) initiative for the improved

530 surveillance of infectious diseases [55, 56]. The AMD and GenomeTrakr

531 frameworks rely on a participatory model with enhanced analysis, curation

532 and data storage at a central site. However, these resource-intensive

533 networks focus on few selected pathogens at present. Notably, there still

534 remain significant challenges for the implementation of the comprehensive

535 WGS services at the local level [48, 57]. It is hoped that the quality

536 framework proposed in the present study would advance the localization of

537 comprehensive WGS services in clinical and public health laboratories.

538      In summary, the salient achievements of this study included: 1)

539 establishment of the performance specifications for WGS in the application to

26

540 public health microbiology in accordance with CLIA guidelines for the LDTs,

541 2) the development of quality assurance (QA) and quality control (QC)

542 measurements for WGS, 3) formatting of laboratory reports for end users

543 with or without WGS expertise, 4) a set of pathogenic bacteria for further

544 validations of WGS and multi-laboratory comparisons and, 5) development

545 of an integrated workflow for the 'wet bench' and 'dry bench' parts of WGS.

546

547 **ACKNOWLEDGEMENT**

555

556 **References**

557 1. Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW: **Transforming clinical microbiology with**
558 **bacterial genome sequencing**. *Nat Rev Genet* 2012, **13**(9):601-612.
559 2. Köser CU, Ellington MJ, Cartwright EJ, Gillespie SH, Brown NM, Farrington M, Holden MT,
560 Dougan G, Bentley SD, Parkhill J: **Routine use of microbial whole genome sequencing in**
561 **diagnostic and public health microbiology**. *PLoS pathogens* 2012, **8**(8):e1002824.
562 3. Grad YH, Lipsitch M, Feldgarden M, Arachchi HM, Cerqueira GC, FitzGerald M, Godfrey P, Haas
563 BJ, Murphy CI, Russ C *et al*: **Genomic epidemiology of the Escherichia coli O104:H4 outbreaks in**
564 **Europe, 2011**. *Proceedings of the National Academy of Sciences* 2012, **109**(8):3065-3070.
565 4. McGann P, Bunin JL, Snesrud E, Singh S, Maybank R, Ong AC, Kwak YI, Seronello S, Clifford RJ,
566 Hinkle M *et al*: **Real time application of whole genome sequencing for outbreak investigation –**
567 **What is an achievable turnaround time?** *Diagnostic Microbiology and Infectious Disease* 2016,
568 **85**(3):277-282.
569 5. Gordon NC, Price JR, Cole K, Everitt R, Morgan M, Finney J, Kearns AM, Pichon B, Young B,
570 Wilson DJ *et al*: **Prediction of Staphylococcus aureus Antimicrobial Resistance by Whole-**
571 **Genome Sequencing**. *Journal of Clinical Microbiology* 2014, **52**(4):1182-1191.
572 6. Leopold SR, Goering RV, Witten A, Harmsen D, Mellmann A: **Bacterial Whole-Genome**
573 **Sequencing Revisited: Portable, Scalable, and Standardized Analysis for Typing and Detection**

574    **of Virulence and Antibiotic Resistance Genes**. *Journal of Clinical Microbiology* 2014, **52**(7):2365-
575    2370.
576  7.   Goodwin S, McPherson JD, McCombie WR: **Coming of age: ten years of next-generation**
577    **sequencing technologies**. *Nature Reviews Genetics* 2016, **17**(6):333-351.
578  8.   Roetzer A, Diel R, Kohl TA, Rückert C, Nübel U, Blom J, Wirth T, Jaenicke S, Schuback S, Rüsch-
579    Gerdes S: **Whole genome sequencing versus traditional genotyping for investigation of a**
580    ***Mycobacterium tuberculosis* outbreak: a longitudinal molecular epidemiological study**. *PLoS*
581    *Med* 2013, **10**(2):e1001387.
582  9.   Gardy JL, Johnston JC, Sui SJH, Cook VJ, Shah L, Brodkin E, Rempel S, Moore R, Zhao Y, Holt
583    R *et al*: **Whole-Genome Sequencing and Social-Network Analysis of a Tuberculosis Outbreak**.
584    *New England Journal of Medicine* 2011, **364**(8):730-739.
585  10.  Eyre DW, Cule ML, Wilson DJ, Griffiths D, Vaughan A, O'Connor L, Ip CLC, Golubchik T, Batty
586    EM, Finney JM *et al*: **Diverse Sources of C. difficile Infection Identified on Whole-Genome**
587    **Sequencing**. *New England Journal of Medicine* 2013, **369**(13):1195-1205.
588  11.  Etienne KA, Roe CC, Smith RM, Vallabhaneni S, Duarte C, Escandón P, Castañeda E, Gómez BL, de
589    Bedout C, López LF *et al*: **Whole-Genome Sequencing to Determine Origin of Multinational**
590    **Outbreak of Sarocladium kiliense Bloodstream Infections**. *Emerging Infectious Diseases* 2016,
591    **22**(3):476-481.
592  12.  Pallen MJ: **Diagnostic metagenomics: potential applications to bacterial, viral and parasitic**
593    **infections**. *Parasitology* 2014, **141**(Special Issue 14):1856-1862.
594  13.  Onderdonk AB, Delaney ML, Fichorova RN: **The Human Microbiome during Bacterial Vaginosis**.
595    *Clinical Microbiology Reviews* 2016, **29**(2):223-238.
596  14.  Pearce MM, Hilt EE, Rosenfeld AB, Zilliox MJ, Thomas-White K, Fok C, Kliethermes S,
597    Schreckenberger PC, Brubaker L, Gai X *et al*: **The Female Urinary Microbiome: a Comparison of**
598    **Women with and without Urgency Urinary Incontinence**. *mBio* 2014, **5**(4).
599  15.  Schubert AM, Rogers MAM, Ring C, Mogle J, Petrosino JP, Young VB, Aronoff DM, Schloss PD:
600    **Microbiome Data Distinguish Patients with Clostridium difficile Infection and Non-C. difficile-**
601    **Associated Diarrhea from Healthy Controls**. *mBio* 2014, **5**(3).
602  16.  Naccache SN, Federman S, Veeraraghavan N, Zaharia M, Lee D, Samayoa E, Bouquet J, Greninger
603    AL, Luk K-C, Enge B *et al*: **A cloud-compatible bioinformatics pipeline for ultrarapid pathogen**
604    **identification from next-generation sequencing of clinical samples**. *Genome Research* 2014,
605    **24**(7):1180-1192.
606  17.  Hasman H, Saputra D, Sicheritz-Ponten T, Lund O, Svendsen CA, Frimodt-Møller N, Aarestrup
607    FM: **Rapid whole genome sequencing for the detection and characterization of**
608    **microorganisms directly from clinical samples**. *Journal of Clinical Microbiology* 2013.
609  18.  Loman NJ, Constantinidou C, Christner M, et al.: **A culture-independent sequence-based**
610    **metagenomics approach to the investigation of an outbreak of shiga-toxigenic escherichia coli**
611    **o104:h4**. *JAMA* 2013, **309**(14):1502-1510.
612  19.  Endrullat C, Glökler J, Franke P, Frohme M: **Standardization and quality management in next-**
613    **generation sequencing**. *Applied & Translational Genomics* 2016.
614  20.  Salto-Tellez M, Gonzalez de Castro D: **Next-generation sequencing: a change of paradigm in**
615    **molecular diagnostic validation**. *J Pathol* 2014, **234**(1):5-10.
616  21.  Goldberg B, Sichtig H, Geyer C, Ledeboer N, Weinstock GM: **Making the Leap from Research**
617    **Laboratory to Clinic: Challenges and Opportunities for Next-Generation Sequencing in**
618    **Infectious Disease Diagnostics**. *mBio* 2015, **6**(6).
619  22.  Kwong JC, McCallum N, Sintchenko V, Howden BP: **Whole genome sequencing in clinical and**
620    **public health microbiology**. *Pathology* 2015, **47**(3):199-210.

621  23.  Luheshi LM, Raza S, Peacock SJ: **Moving pathogen genomics out of the lab and into the clinic:**
622       **what will it take?** *Genome Medicine* 2015, **7**(1):1-3.
623  24.  Oliver GR, Hart SN, Klee EW: **Bioinformatics for Clinical Next Generation Sequencing**. *Clinical*
624       *Chemistry* 2015, **61**(1):124-135.
625  25.  Fricke WF, Rasko DA: **Bacterial genome sequencing in the clinic: bioinformatic challenges and**
626       **solutions**. *Nat Rev Genet* 2014, **15**(1):49-55.
627  26.  Wyres KL, Conway TC, Garg S, Queiroz C, Reumann M, Holt K, Rusu LI: **WGS Analysis and**
628       **Interpretation in Clinical and Public Health Microbiology Laboratories: What Are the**
629       **Requirements and How Do Existing Tools Compare?**
630  27.  Rhoads DD, Sintchenko V, Rauch CA, Pantanowitz L: **Clinical Microbiology Informatics**. *Clinical*
631       *Microbiology Reviews* 2014, **27**(4):1025-1047.
632  28.  Gargis AS, Kalman L, Lubin IM: **Assuring the Quality of Next-Generation Sequencing in Clinical**
633       **Microbiology and Public Health Laboratories**. *Journal of Clinical Microbiology* 2016:JCM. 00949-
634       00916.
635  29.  Lefterova MI, Suarez CJ, Banaei N, Pinsky BA: **Next-Generation Sequencing for Infectious**
636       **Disease Diagnosis and Management: A Report of the Association for Molecular Pathology**. *The*
637       *Journal of molecular diagnostics : JMD* 2015, **17**(6):623-634.
638  30.  Aziz N, Zhao Q, Bry L, Driscoll DK, Funke B, Gibson JS, Grody WW, Hegde MR, Hoeltge GA,
639       Leonard DGB *et al*: **College of American Pathologists' Laboratory Standards for Next-**
640       **Generation Sequencing Clinical Tests**. *Archives of pathology & laboratory medicine* 2014,
641       **139**(4):481-493.
642  31.  Moran-Gilad J, Sintchenko V, Pedersen SK, Wolfgang WJ, Pettengill J, Strain E, Hendriksen RS:
643       **Proficiency testing for bacterial whole genome sequencing: an end-user survey of current**
644       **capabilities, requirements and priorities**. *BMC Infectious Diseases* 2015, **15**(1):1-10.
645  32.  Olson ND, Jackson SA, Lin NJ: **Report from the Standards for Pathogen Identification via Next-**
646       **Generation Sequencing (SPIN) Workshop**. *Standards in Genomic Sciences* 2015, **10**(119).
647  33.  Gargis AS, Kalman L, Berry MW, Bick DP, Dimmock DP, Hambuch T, Lu F, Lyon E, Voelkerding KV,
648       Zehnbauer BA: **Assuring the quality of next-generation sequencing in clinical laboratory**
649       **practice**. *Nature biotechnology* 2012, **30**(11):1033-1036.
650  34.  Gargis AS, Kalman L, Bick DP, da Silva C, Dimmock DP, Funke BH, Gowrisankar S, Hegde MR,
651       Kulkarni S, Mason CE *et al*: **Good laboratory practice for clinical next-generation sequencing**
652       **informatics pipelines**. *Nature biotechnology* 2015, **33**(7):689-693.
653  35.  Duncavage EJ, Abel HJ, Merker JD, Bodner JB, Zhao Q, Voelkerding KV, Pfeifer JD: **A Model Study**
654       **of In Silico Proficiency Testing for Clinical Next-Generation Sequencing**. *Archives of pathology &*
655       *laboratory medicine* 2016, **140**(10):1085-1091.
656  36.  CLSI: **Nucleic Acid Sequencing Methods in Diagnostic Laboratory Medicine: Approved**
657       **Guideline- 2d edition. MM09-A2**. 2014.
658  37.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
659       Genome Project Data Processing S: **The Sequence Alignment/Map format and SAMtools**.
660       *Bioinformatics* 2009, **25**(16):2078-2079.
661  38.  Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G,
662       Marth GT, Sherry ST *et al*: **The variant call format and VCFtools**. *Bioinformatics* 2011,
663       **27**(15):2156-2158.
664  39.  Seemann T: **Prokka: rapid prokaryotic genome annotation**. *Bioinformatics* 2014, **30**(14):2068-
665       2069.
666  40.  Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, Brown CT, Porras-Alfaro A, Kuske CR,
667       Tiedje JM: **Ribosomal Database Project: data and tools for high throughput rRNA analysis**.
668       *Nucleic acids research* 2014, **42**(Database issue):D633-642.

669 41.   Larsen MV, Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, Jelsbak L, Sicheritz-Ponten
670       T, Ussery DW, Aarestrup FM *et al*: **Multilocus sequence typing of total-genome-sequenced**
671       **bacteria**. *J Clin Microbiol* 2012, **50**(4):1355-1361.
672 42.   Joensen KG, Scheutz F, Lund O, Hasman H, Kaas RS, Nielsen EM, Aarestrup FM: **Real-time whole-**
673       **genome sequencing for routine typing, surveillance, and outbreak detection of verotoxigenic**
674       **Escherichia coli**. *J Clin Microbiol* 2014, **52**(5):1501-1510.
675 43.   CLSI: **Performance Standards for Antimicrobial Susceptibility Testing; Twenty-Fifth**
676       **Informational Supplement. CLSI document M100-S25. Wayne, PA: CLinical and Laboratory**
677       **STandards Institute; 2015**. 2015.
678 44.   Harris SR, Cartwright EJ, Torok ME, Holden MT, Brown NM, Ogilvy-Stuart AL, Ellington MJ, Quail
679       MA, Bentley SD, Parkhill J *et al*: **Whole-genome sequencing for analysis of an outbreak of**
680       **meticillin-resistant Staphylococcus aureus: a descriptive study**. *The Lancet Infectious diseases*
681       2013, **13**(2):130-136.
682 45.   Leekitcharoenphon P, Nielsen EM, Kaas RS, Lund O, Aarestrup FM: **Evaluation of whole genome**
683       **sequencing for outbreak detection of Salmonella enterica**. *PloS one* 2014, **9**(2):e87991.
684 46.   CLSI: **Molecular Methods for Bacterial Strain Typing; Approved Guideline, MM11-A**. 2007.
685 47.   Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, Grody WW, Hegde M, Lyon E, Spector E
686       *et al*: **Standards and guidelines for the interpretation of sequence variants: a joint consensus**
687       **recommendation of the American College of Medical Genetics and Genomics and the**
688       **Association for Molecular Pathology**. *Genet Med* 2015, **17**(5):405-423.
689 48.   Lesho E, Clifford R, Onmus-Leone F, Appalla L, Snesrud E, Kwak Y, Ong A, Maybank R, Waterman
690       P, Rohrbeck P *et al*: **The Challenges of Implementing Next Generation Sequencing Across a**
691       **Large Healthcare System, and the Molecular Epidemiology and Antibiotic Susceptibilies of**
692       **Carbapenemase-Producing Bacteria in the Healthcare System of the U.S. Department of**
693       **Defense**. *PloS one* 2016, **11**(5):e0155770.
694 49.   Kalman LV, Datta V, Williams M, Zook JM, Salit ML, Han J-Y: **Development and Characterization**
695       **of Reference Materials for Genetic Testing: Focus on Public Partnerships**. *Annals of Laboratory*
696       *Medicine* 2016, **36**(6):513-520.
697 50.   Lindsey RL, Pouseele H, Chen JC, Strockbine NA, Carleton HA: **Implementation of Whole**
698       **Genome Sequencing (WGS) for Identification and Characterization of Shiga Toxin-Producing**
699       **Escherichia coli (STEC) in the United States**. *Frontiers in Microbiology* 2016, **7**(766).
700 51.   Pankhurst LJ, del Ojo Elias C, Votintseva AA, Walker TM, Cole K, Davies J, Fermont JM, Gascoyne-
701       Binzi DM, Kohl TA, Kong C *et al*: **Rapid, comprehensive, and affordable mycobacterial diagnosis**
702       **with whole-genome sequencing: a prospective study**. *The Lancet Respiratory Medicine* 2016,
703       **4**(1):49-58.
704 52.   Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, Iqbal Z, Feuerriegel S, Niehaus
705       KE, Wilson DJ *et al*: **Whole-genome sequencing for prediction of *Mycobacterium tuberculosis***
706       **drug susceptibility and resistance: a retrospective cohort study**. *The Lancet Infectious Diseases*
707       2015, **15**(10):1193-1202.
708 53.   Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, Aarestrup FM, Larsen
709       MV: **Identification of acquired antimicrobial resistance genes**. *The Journal of antimicrobial*
710       *chemotherapy* 2012, **67**(11):2640-2644.
711 54.   Arnold C: **Considerations in centralizing whole genome sequencing for microbiology in a public**
712       **health setting**. *Expert Review of Molecular Diagnostics* 2016, **16**(6):619-621.
713 55.   Auffray C, Caulfield T, Griffin JL, Khoury MJ, Lupski JR, Schwab M: **From genomic medicine to**
714       **precision medicine: highlights of 2015**. *Genome Medicine* 2016, **8**(1):12.

715  56.  Allard MW, Strain E, Melka D, Bunning K, Musser SM, Brown EW, Timme R: **Practical Value of Food Pathogen Traceability through Building a Whole-Genome Sequencing Network and Database**. *Journal of Clinical Microbiology* 2016, **54**(8):1975-1983.

718  57.  Robilotti E, Kamboj M: **Integration of Whole-Genome Sequencing into Infection Control Practices: the Potential and the Hurdles**. *Journal of Clinical Microbiology* 2015, **53**(4):1054-1055.

721

722

723    **FIGURE LEGENDS**

724
725    Figure 1. **WGS wet and dry bench workflow**
726

727    Figure 2. **The summary of the WGS validation.**

728    Figure 3. **Bioinformatics pipeline validation with two groups of**

729    **outbreak isolates. A.** "Study 1 tree", a phylogenetic tree of outbreak

730    isolates, which was published in the study 1. The isolates from the study

731    which were picked for validation have arrows pointing at them and numbers

732    assigned for purposes of validation (1-7). **B.** A tree representing

733    phylogenetic connections between chosen isolates from original study tree.

734    **C.** "Validation tree 1", a phylogenetic tree generated using the PHL

735    bioinformatics pipeline. The same isolates in the original tree and validation

736    tree are marked with the same numbers. **D.** Comparison of the group of

737    related isolates (1-7) from Study 1 with epidemiologically unrelated isolates

738    from the same study using the PHL bioinformatics pipeline. **E.** "Study two

739    tree", a phylogenetic tree combining epidemiologically related and

740    nonrelated isolates published in the study 2. The isolates from the study two

741    which were picked for validation marked with green node circles and had

742    numbers 1-11 assigned for purposes of validation. **F.** A tree representing

743    phylogenetic connections between chosen isolates from original study tree.

744    **G.** "Validation tree 2", a phylogenetic tree generated using the PHL

745    bioinformatics pipeline. The same isolates in the tree from Study 2 and the

746    validation tree are marked with the same numbers.

32

747    Supplementary Figure 1. **WGS quality control scheme.**

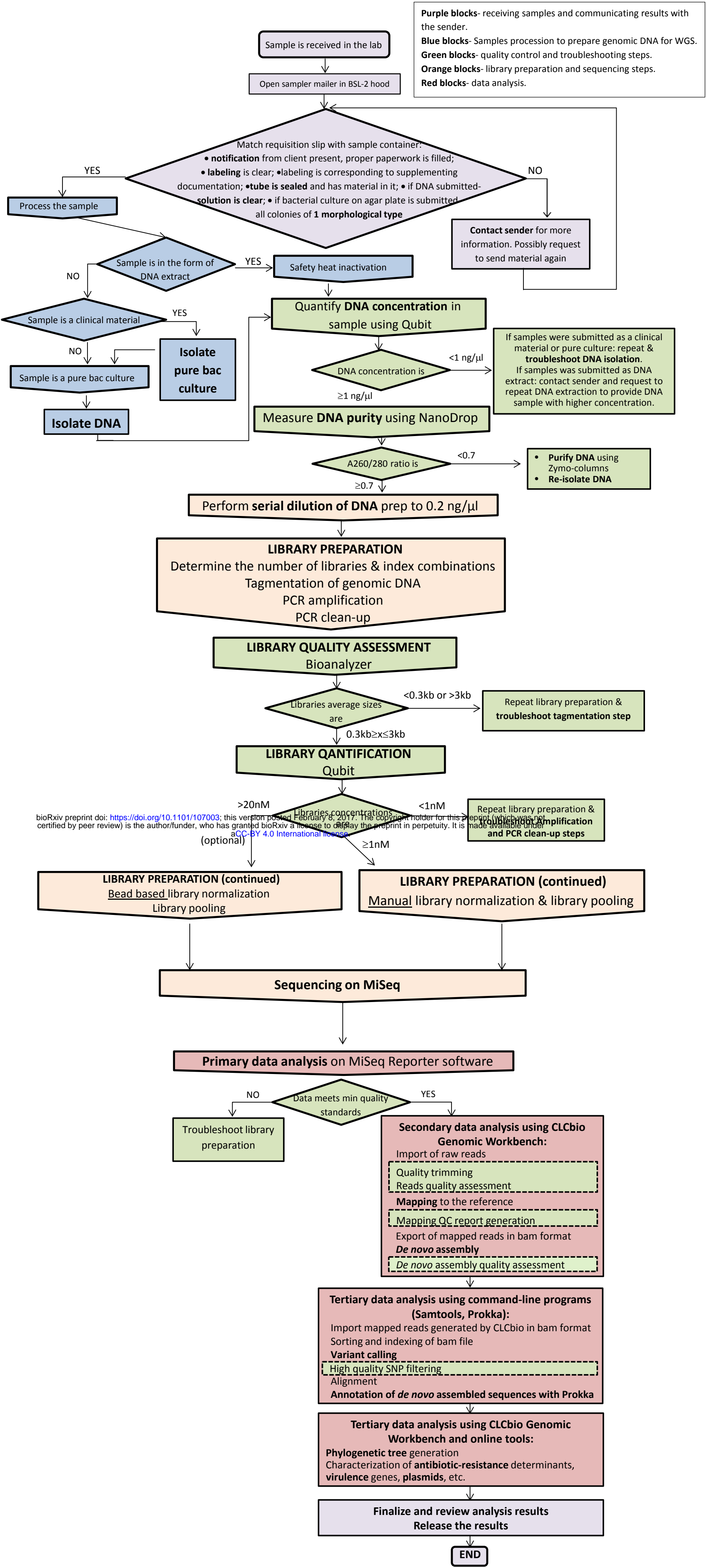748

749

Figure 1. **WGS wet and dry bench workflow**



**Purple blocks**- receiving samples and communicating results with the sender.
**Blue blocks**- Samples procession to prepare genomic DNA for WGS.
**Green blocks**- quality control and troubleshooting steps.
**Orange blocks**- library preparation and sequencing steps.
**Red blocks**- data analysis.

Sample is received in the lab

Open sampler mailer in BSL-2 hood

Match requisition slip with sample container:
• **notification** from client present, proper paperwork is filled;
• **labeling** is clear; •labeling is corresponding to supplementing documentation; •**tube is sealed** and has material in it; • if DNA submitted- **solution is clear**; • if bacterial culture on agar plate is submitted all colonies of **1 morphological type**

YES

NO

Process the sample

Contact sender for more information. Possibly request to send material again

Sample is in the form of DNA extract

YES

Safety heat inactivation

NO

Sample is a clinical material

YES

Isolate pure bac culture

NO

Sample is a pure bac culture

Isolate DNA

Quantify **DNA concentration** in sample using Qubit

DNA concentration is

<1 ng/μl

≥1 ng/μl

If samples were submitted as a clinical material or pure culture: repeat & **troubleshoot DNA isolation**. If samples was submitted as DNA extract: contact sender and request to repeat DNA extraction to provide DNA sample with higher concentration.

Measure **DNA purity** using NanoDrop

A260/280 ratio is

<0.7

• **Purify DNA** using Zymo-columns
• **Re-isolate DNA**

≥0.7

Perform **serial dilution of DNA** prep to 0.2 ng/μl

**LIBRARY PREPARATION**
Determine the number of libraries & index combinations
Tagmentation of genomic DNA
PCR amplification
PCR clean-up

**LIBRARY QUALITY ASSESSMENT**
Bioanalyzer

Libraries average sizes are

<0.3kb or >3kb

Repeat library preparation & **troubleshoot tagmentation step**

0.3kb≥x≤3kb

**LIBRARY QANTIFICATION**
Qubit

Libraries concentrations

>20nM

<1nM

Repeat library preparation & **troubleshoot PCR amplification and PCR clean-up steps**

≥1nM

(optional)

**LIBRARY PREPARATION (continued)**
Bead based library normalization
Library pooling

**LIBRARY PREPARATION (continued)**
Manual library normalization & library pooling

**Sequencing on MiSeq**

**Primary data analysis** on MiSeq Reporter software

Data meets min quality standards

NO

Troubleshoot library preparation

YES

**Secondary data analysis using CLCbio Genomic Workbench:**
Import of raw reads
Quality trimming
Reads quality assessment
**Mapping** to the reference
Mapping QC report generation
Export of mapped reads in bam format
*De novo* assembly
*De novo* assembly quality assessment

**Tertiary data analysis using command-line programs (Samtools, Prokka):**
Import mapped reads generated by CLCbio in bam format
Sorting and indexing of bam file
**Variant calling**
High quality SNP filtering
Alignment
**Annotation of *de novo*** assembled sequences with Prokka

**Tertiary data analysis using CLCbio Genomic Workbench and online tools:**
**Phylogenetic tree** generation
Characterization of **antibiotic-resistance** determinants, **virulence** genes, **plasmids**, etc.

**Finalize and review analysis results**
**Release the results**

END

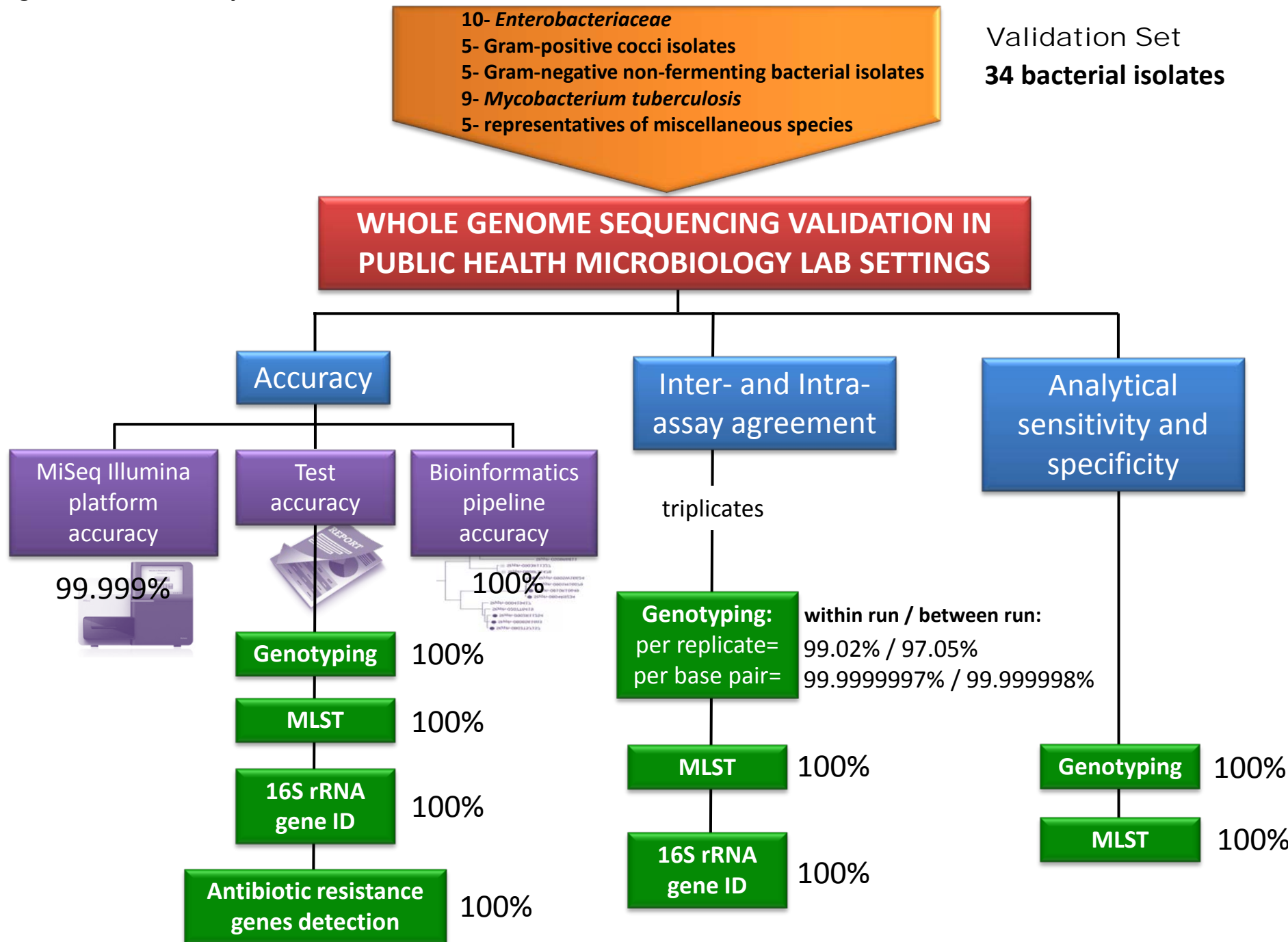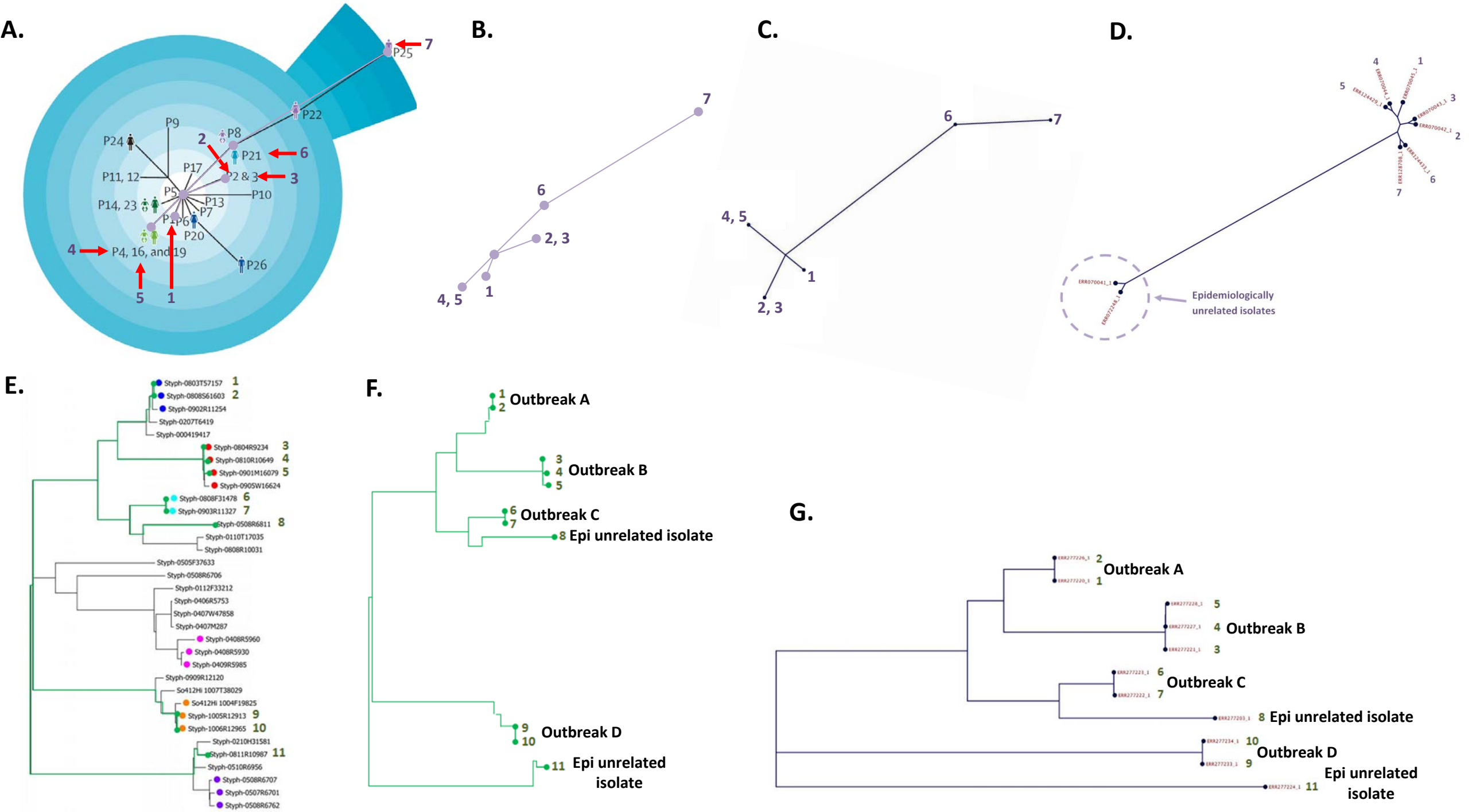Figure 2. **The summary of the WGS validation.**

Figure 3. **Bioinformatics pipeline validation with two groups of outbreak isolates.**

**Table 1. List of strains used for validation and corresponding reference materials**

| MDL ID | Species | Reference | |
|---|---|---|---|
| | | **NCBI Strain** | **NCBI Acc#** |
| | | *Reference materials- NCBI strains* | |
| C1 | *Escherichia coli* O157:H7 CDC EDL 933 | O157:H7 CDC EDL 933 | NZ_CP008957.1 |
| C3 | *Escherichia coli* ATCC 8739 | ATCC 8739 | NC_010468.1 |
| C55 | *Escherichia coli* ATCC 25922 | ATCC 25922 | NZ_CP009072.1 |
| C4 | *Enterobacter cloacae* ATCC 13047 | ATCC 13047 | NC_014121 |
| C6 | *Salmonella enterica* ser Typhimurium ATCC 14028 | 14028S | NC_016856 |
| C5 | *Staphylococcus aureus* ATCC 25923 | ATCC 25923 | NZ_CP009361 |
| C46 | *Enterococcus faecalis* ATCC 29212 | ATCC 29212 | NZ_CP008816 |
| C47 | *Staphylococcus epidermidis* ATCC 12228 | ATCC 12228 | NC_004461 |
| C48 | *Staphylococcus saprophyticus* ATCC 15305 | ATCC 15305 | NC_007350 |
| C49 | *Streptococcus pneumoniae* ATCC 6305 | ATCC 700669 | FM211187 |
| C50 | *Pseudomonas aeruginosa* ATCC 27853 | FRD1 | NZ_CP010555 |
| C51 | *Stenotrophomonas maltophilia* ATCC 13637 | ATCC 13637 | NZ_CP008838 |
| C52 | *Legionella pneumophila* SG-12 ATCC 43290 | ATCC 43290 | NC_016811 |
| C53 | *Moraxella catarrhalis* 87A-3084 | ATCC 25240 | NZ_CP008804 |
| C54 | *Acinetobacter baumannii* ATCC 17945 | AB07 | NZ_CP006963 |
| C103 | *Bacteroides fragilis* ATCC 25285 | 638R | NC_016776 |
| C104 | *Haemophilus influenzae* ATCC 10211 | KR494 | NC_022356 |
| C2 | *Aeromonas hydrophilia* ATCC 7966 | ATCC 7966 | NC_008570 |
| C105 | *Corynebacterium jeikeium* ATCC 43734 | ATCC 43734 | GG700813:GG700833 |
| C106 | *Neisseria gonorrhoeae* ATCC 49226 | MS11 | NC_022240 |
| C56 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C57 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C58 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C59 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C61 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C65 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C67 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C68 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |
| C69 | *Mycobacterium tuberculosis* | H37Rv | NC_000962.3 |

| MDL ID | Species | Reference raw reads generated by CDC | | Reference used for mapping | |
|---|---|---|---|---|---|
| | | **CDC Strain** | **Accession #** | **NCBI Strain** | **NCBI Acc#** |
| | | *Reference materials- strains sequenced at CDC* | | | |
| C72 | *Escherichia coli* O121:H19 | 2014C-3857 | SRR1610033 | 2011C-3493 | NC_018658 |
| C73 | *Salmonella enterica* ser Enteritidis | CDC_2010K-1543 | SRR518749 | P125109 | NC_011294.1 |
| C74 | *Salmonella enterica* ser Infantis | 2014K-0434 | SRR1616809 | 1326/28 | NZ_LN649235 |
| C75 | *Salmonella enterica* ser Adelaide | 2014K-0941 | SRR1686419 | P125109 | NC_011294.1 |
| C76 | *Salmonella enterica* ser Worthington | 2012K-1219 | SRR1614868 | P125109 | NC_011294.1 |
| C77* | *Salmonella enterica* ser Saintpaul | 2014K-0875 | SRR1640105 | 14028S | NC_016856 |

**Footnotes**: Green color designates cases when the genome of the strain sequenced by the PHL is available from the NCBI database. Yellow color designates cases when the genome of the strain sequenced by the PHL is NOT available from the NCBI database and an alternative reference genome was used for mapping.

*P.S.: Sample C77 was sequenced by PHL only for genotyping test accuracy validation. No replicates were done.

**Table 2. Performance characteristics, definitions, and formulas used in validation. Summary of the validation for different assays.**

| | | | Definition of the performance characteristic for WGS applications | Formula used for calculation | Assay used for validation of the parameter | | | | Assay-specific definitions |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | hqSNP based genotyping | MLST | 16S | Antibiotic resistance genes detection | |
| Accuracy | The degree of agreement between the nucleic acid sequences derived from the assay (measured value) and a reference sequence (the true value). | accuracy of platform | Accuracy of base calling against reference sequence. The accuracy of the platform was established by determining the proximity of agreement between base calling made by MiSeq sequencer (measured value) and NCBI reference sequence (the true value). | % agreement with reference = $\dfrac{\text{(Covered genome length)} - \text{(Total \# of SNP differing from reference)}}{\text{Covered genome length}}$ x 100% | 99.999378% | | | | Accuracy of the platform |
| | | accuracy of assay | Test accuracy is determined by agreement of the test result of validation samples sequenced by PHL with the test result for reference sequences of the same strains. | $\text{Accuracy} = \dfrac{\text{\# of correct results}}{\text{total \# of results}}$ x 100% | congruence of phylogenetic trees built using reference sequences and validation sequences | Detection and correct identification of each of the MLST alleles | ID of 16S rRNA sequence of the validation sample matches the ID of 16S rRNA sequence of the reference sequence | Presence of ABR genes characteristic for reference strain, absence of any other ABR genes | Definition of correct results |
| | | | | | Individual sample clustering | Allele | 16S rRNA ID result | Antibiotic resistance gene | Single test unit |
| | | | | | 100% | 100% | 100% | 100% | Accuracy of the assay |
| | | accuracy of bioinformatics pipeline | Clustering suggested by previous investigators must match clustering achieved by the analysis using PHL validation bioinformatics pipeline. | % agreement = $\dfrac{\text{\# of outbreak isolates clustered correctly in validation tree}}{\text{Total \# of outbreak isolates clustered together in the study tree}}$ x 100% | 100% | - | - | - | Accuracy of the bioinformatics pipeline |
| Precision | The degree to which repeated sequence analyses give the same result repeatability (within-run precision) and reproducibility (between-run precision). | Repeatability (precision within run) | was established by sequencing the same samples multiple times under the same conditions and evaluating the concordance of the test results and performance. | $\text{Inter} - \text{assay precision (Repeatability)} =$ $\dfrac{\text{\# within-run replicates in agreement}}{\text{Total \# of tests performed for within-run replicates}}$ x 100% | inter- assay precision of single nucleotide variant detection. | | | | Definition of correct results |
| | | | | | SNP (precision per replicate) | SNP (precision per base pair) | Allele | 16S rRNA ID | - | Single test unit |
| | | | | | 99.02% | 99.9999997% | 100% | 100% | - | Repeatability |
| | | Reproducibility (precision between runs) | was assessed as the consistency of the test results and performance characteristics for the same sample sequenced under different conditions, such as between different runs and different sample preparations. | $\text{Intra} - \text{assay precision (Reproducibility)} =$ $\dfrac{\text{\# between-run replicates in agreement}}{\text{Total \# of tests performed for between-run replicates}}$ x 100% | intra- assay precision of single nucleotide variant detection. | | | | Definition of correct results |
| | | | | | SNP (precision per replicate) | SNP (precision per base pair) | Allele | 16S rRNA ID | - | Single test unit |
| | | | | | 97.05% | 99.999998% | 100% | 100% | - | Reproducibility |

Continued

| | Definition of the performance characteristic for WGS applications | Formula used for calculation | Assay used for validation of the parameter | | | | Assay-specific definitions |
|---|---|---|---|---|---|---|---|
| | | | hqSNP based genotyping | MLST | 16S | Antibiotic resistance genes detection | |
| Analytical sensitivity (Limit of detection) | The likelihood that a WGS assay will detect a sequence variation when present within the analyzed genomic region (this value reflects a false negative rate of the test). | $\text{Analytical sensitivity} = \dfrac{TP}{TP + FN} \times 100\%$ | Clustering of related samples (#of validation samples with clustering results matching reference ) | Number of correctly identified alleles | - | - | Definition of true positive results |
| | | | Number of validation samples which clustered together with samples, genetically distant according to the reference tree | Number of unidentified or misidentified alleles validation samples | - | - | Definition of false negative results |
| | | | Individual sample clustering | Allele | - | - | Single test unit |
| | | | 100% | 100% | - | - | Analytical sensitivity |
| Analytical specificity | The probability that a WGS assay will not detect sequence variation(s) when none are present within the analyzed genomic region (this value reflects a test's false positive rate). | $\text{Analytical specificity} = \dfrac{TN}{TN + FP} \times 100\%$ | No clustering between unrelated samples (#of validation samples with clustering results matching reference ) | Number of unidentified alleles in negative control samples | - | - | Definition of true negative results |
| | | | Number of validation samples which failed to clustered together with samples, genetically similar according to the reference tree | Number of identified alleles in negative control samples | - | - | Definition of false positive results |
| | | | Individual sample clustering | Allele | - | - | Single test unit |
| | | | 100% | 100% | - | - | Analytical specificity |
| Reportable range | The region of the genome in which sequence of an acceptable quality can be derived by the laboratory test. | N/A | Genome-wide hq SNPs | Housekeeping genes used in corresponding MLST schemes | 16S rRNA gene | Antibiotic resistance genes in included ResFinder database | |

**Footnotes**: See details in Supplementary Document 1.

**Abbreviations**: TP- True positive results, TN- True negative results, FP- False positive, FN- False negative

## Table 3. Summary of the studies used for validation of bioinformatics pipeline

| Study | Study 1.<br>SR Harris et al. Lancet Infect Dis 2013; 13: 130–36 [44] | Study 2.<br>P Leekitcharoenphon et al. PLoS ONE 2014; 9(2): e87991 [45] |
|---|---|---|
| **Microorganism** | Methicillin-resistant *Staphylococcus aureus* | *Salmonella enterica* serovar Typhimurium |
| **Source of isolates** | Human | Human |
| **Number of isolates analyzed** | 7 outbreak isolates (1 outbreak cluster) + 2 epidemiologically unrelated isolates | 9 outbreak isolates (4 outbreak clusters) + 2 epidemiologically unrelated isolates |
| **Type of outbreak** | Hospital-associated outbreak | Foodborne outbreaks |
| **ID of the samples in the study which were used for validation** | P1, P2, P3, P4, P16, P21, P25, Identified by Infectious Control Investigation non-outbreak ST1, MRSA identified by searching microbiology database non-outbreak ST772 | 0803T57157, 0808S61603, 0808F31478, 0903R11327, 0811R10987, 0804R9234, 0810R10649, 0901M16079, 0110T17035, 1005R12913, 1006R12965 |
| **Accession ## of corresponding samples** | ERR070045, ERR070042, ERR070043, ERR070044, ERR124429, ERR124433, ERR128708, ERR070041, ERR072248 | ERR277220, ERR277226, ERR277223, ERR277222, ERR277224, ERR277221, ERR277227, ERR277228, ERR277203, ERR277233, ERR277234 |
| **# of clusters in study tree** | 1 | 4 |
| **# of clusters in validation tree** | 1 | 4 |
| **# of outbreak isolates in each cluster in the study tree** | Cluster 1= 7 | Cluster 1= 2, Cluster 2= 3, Cluster 3= 2, Cluster 4= 2 |
| **# of outbreak isolates in each cluster in validation tree** | Cluster 1= 7 | Cluster 1= 2, Cluster 2= 3, Cluster 3= 2, Cluster 4= 2 |
| **# of epidemiologically unrelated isolates in the set** | 2 | 2 |
| **# of epidemiologically unrelated isolates clustered with outbreak isolates** | 0 | 0 |
| **% agreement= (# of outbreak isolates clustered correctly in validation tree)x100%/ (Total # of outbreak isolates clustered together in the study tree)** | (7x100/7) = 100% | (9x100/9) = 100% |

# WGS QUALITY CONTROL SCHEME

| | Each tested samples | Positive control | | Negative control | |
|---|---|---|---|---|---|

**DNA template QC**

| | | *E. coli ATCC 25922* | *Listeria H8394 or M.tuberculosis [DNA extraction control]* | | *No-template control* |

- DNA concentration
- DNA purity

**Library QC**

- Library size distribution
- Library concentration

**Run QC**

- Percent of bases with quality score >Q30
- Cluster density
- Cluster passing filter

**Raw Data QC**

| Each tested samples | Positive control | Negative control | |
|---|---|---|---|
| •Average depth of the genome coverage<br>•Accuracy of base calling: Read length with quality score ≥Q30 | *Spiked-in positive PhiX control*<br>•Error rate | *Negative index combination control*<br>•Number of reads after trim | •Number of reads after trim |

**Analysis QC**

| Each tested samples | Positive control | Negative control | |
|---|---|---|---|
| **16S  rRNA ID**  [match submitter ID if available] | | •N50 for *de novo* assembled reads<br>•The highest coverage of *de novo* assembled contigs | •N50 for *de novo* assembled reads<br>•The highest coverage of *de novo* assembled contigs |
| ***De novo* QC:**<br>• Min contig coverage<br>• Min contig length | | | |
| **Mapping QC:**<br>• Uniformity of coverage<br>• Percentage of genome covered | •MLST allelic profile<br>•Antibiotic resistance genes<br>•Virulence genes | | |
| **Variant calling QC:**<br>•Min SNP coverage<br>•Min SNP quality<br>•No heterozygotes | | *Negative control for phylogenetic analysis*<br>**Epidemiologically unrelated control strain of bacteria** | |

## Workflow (left column)

- Bacterial culture
- DNA isolation
- Library preparation
- Sequencing on MiSeq
- Primary data analysis
- Secondary data analysis
  - Raw reads import Trimming
    - •Mapping to reference
    - •hqSNP calling
    - •Phylogenetic tree generation
    - •*De novo* assembly
    - •Gene annotation
    - •16S rRNA ID
    - •*In silico* MLST
    - •ABR/Virulence genes detection
- Report generation