

1 **Title: FDA-ARGOS: A Public Quality-Controlled Genome Database Resource for Infectious**  
2 **Disease Sequencing Diagnostics and Regulatory Science Research**

3  
4 Heike Sichtig<sup>1\*</sup>, Timothy Minogue<sup>2\*</sup>, Yi Yan<sup>1</sup>, Christopher Stefan<sup>2</sup>, Adrienne Hall<sup>2</sup>, Luke Tallon<sup>3</sup>,  
5 Lisa Sadzewicz<sup>3</sup>, Suvarna Nadendla<sup>3</sup>, William Klimke<sup>4</sup>, Eneida Hatcher<sup>4</sup>, Martin Shumway<sup>4</sup>,  
6 Dayanara Lebron Aldea<sup>5</sup>, Jonathan Allen<sup>5</sup>, Jeffrey Koehler<sup>2</sup>, Tom Slezak<sup>5</sup>, Stephen Lovell<sup>1</sup>, Randal  
7 Schoepp<sup>2</sup> and Uwe Scherf<sup>1</sup>

8 <sup>1</sup> U.S. Food and Drug Administration, <sup>2</sup> U.S. Army Medical Research Institute of Infectious Diseases, <sup>3</sup>  
9 Institute for Genome Sciences at the University of Maryland, <sup>4</sup> National Center for Biotechnology  
10 Information, National Library of Medicine, National Institutes of Health, <sup>5</sup> Lawrence Livermore National  
11 Laboratories

12 \*Correspondence: [Heike.Sichtig@fda.hhs.gov](mailto:Heike.Sichtig@fda.hhs.gov), [Timothy.D.Minogue.civ@mail.mil](mailto:Timothy.D.Minogue.civ@mail.mil)

13  
14 **ACCESSION NUMBERS**

15 FDA-ARGOS raw data, assemblies, annotations, metadata and pipeline information are available from  
16 Bioproject ID# PRJNA231221 and at <https://www.ncbi.nlm.nih.gov/bioproject/231221>. Reference data  
17 sets from the use cases are available from Bioproject ID# PRJNA495928 and at  
18 <https://www.ncbi.nlm.nih.gov/bioproject/495928>.

19  
20 **SUPPLEMENTAL INFORMATION**

21 Supplemental Information includes 8 Supplementary Tables, 5 Reference Data Sets and FDA-ARGOS  
22 Wanted Organism List.

23  
24 **AUTHOR CONTRIBUTIONS**

25 H.S. conceived of the project, led the project, collected samples, registered samples, wrote and revised the  
26 manuscript, generated figures and tables, performed data analysis, and served as the principal investigator. T.D.M.  
27 led the coordination of the use cases and wrote and revised the manuscript. Y.Y. did script/command  
28 development, data analysis, gathered and organized FDA-ARGOS database metrics. A.H. performed DNA  
29 extraction, library preparation, MIPS and Illumina sequencing, and gathered and organized data for the Ebola *in*  
30 *silico* study, C.S. collected and isolated samples, extracted DNA, performed library preparations and Illumina  
31 sequencing, gathered and organized data for the *E. avium* gap study, and generated and revised figures, L.T. did  
32 IGS sequencing work, L.S. did IGS sequencing work, S.N. did data analysis, and gathered, registered in NCBI and  
33 organized data from IGS sequencing work, W.K. and M.S. helped with coordination of BioProject and data  
34 submissions and bacterial annotations and the assessment for gap filling, E.H. did viral annotations, D.L. did LMAT  
35 analysis, J.A. coordinated LMAT analysis, J.K. collected and sequenced samples, T.S. helped develop the study and

36 experimental design, S.L. helped develop the study and experimental design, R.S. collected clinical samples, U.S.  
37 helped develop the study and experimental design.  
38

39 **ABSTRACT:**

40 Infectious disease next generation sequencing (ID-NGS) diagnostics are on the cusp of  
41 revolutionizing the clinical market. To facilitate this transition, FDA proactively invested in tools  
42 to support innovation of emerging technologies. FDA and collaborators established a publicly  
43 available database, FDA dAtabase for Regulatory-Grade micrObial Sequences (FDA-ARGOS), as a  
44 tool to fill reference database gaps with quality-controlled genomes. This manuscript discusses  
45 quality control metrics for the proposed FDA-ARGOS genomic resource and outlines the need  
46 for quality-controlled genome gap filling in the public domain. Here, we also present three case  
47 studies showcasing potential applications for FDA-ARGOS in infectious disease diagnostics,  
48 specifically: assay design, reference database and *in silico* sequence comparison in combination  
49 with representative microbial organism wet lab testing; a novel composite validation strategy  
50 for ID-NGS diagnostics. The use of FDA-ARGOS as an *in silico* comparator tool could reduce the  
51 burden for completing ID-NGS clinical trials. In addition, use cases identifying *Enterococcus*  
52 *avium* and Ebola virus (Zaire ebolavirus variant Makona) demonstrate the utility of FDA-ARGOS  
53 as a reference database for independent performance validation of new tests and for  
54 documenting how one would use this database as an *in silico* sequence target comparator tool  
55 for ID-NGS validation, respectively.

56

57 **Key Words:** infectious disease (ID), next generation sequencing (NGS), diagnostics (Dx), ID-NGS-  
58 Dx, agnostic (unbiased) ID sequencing, targeted ID sequencing, metagenomics shotgun  
59 sequencing, isolate shotgun sequencing, *Enterococcus avium*, Ebolavirus

60

61

## 62 **INTRODUCTION:**

63 The Food and Drug Administration's (FDA) premarket review of in vitro diagnostics relies  
64 on safety, efficacy, quality, and performance and ensures patient access to safe and accurate  
65 new technologies, such as next-generation sequencing. Within this premarket review, FDA  
66 performs risk-based evaluation of novel diagnostic devices by leveraging clinical expertise, as  
67 well as research evidence to support regulatory decisions and considers patient values and  
68 preferences. Infectious disease next generation sequencing (ID-NGS) diagnostics, with the  
69 potential to identify any microbial organism or genomic marker from a patient sample in a  
70 single test, are poised to enter the clinical diagnostic laboratory (Goldberg, Sichtig et al. 2015,  
71 Arnold 2017, Heger 2018). For accurate identification of any infectious organism, ID-NGS  
72 requires comprehensive reference databases, thereby strongly emphasizing the need for more  
73 complete high-quality reference genomes. Metagenomic agnostic sequencing also requires  
74 novel validation strategies as the traditional diagnostic evaluation for all known organisms is  
75 unfeasible. Described in greater detail throughout this paper (Figure 1), here we are describing

76 one effort for defining a composite-reference method approach for ID-NGS device validation  
77 utilizing *in silico* sequence comparison.

78 Patients and clinicians need alternative solutions when conventional diagnostics (e.g.,  
79 real-time PCR, culture or ELISA), fail to identify an infectious etiology. Several studies document  
80 this need of applying hypothesis-free NGS as a diagnostic of last resort, such as high-risk  
81 transplant population or failure of diagnosis with conventional diagnostics (Schlaberg, Chiu et  
82 al. 2017, Wilson, Zimmermann et al. 2017). Numerous groups have successfully applied ID-NGS  
83 technology across several unique and diverse clinical use cases. For example, isolate shotgun  
84 sequencing information uncovered unexpected transmission routes during multi-drug resistant  
85 nosocomial organism outbreaks (Snitkin, Zelazny et al. 2012, Roach, Burton et al. 2015, Snitkin,  
86 Won et al. 2017). Other studies showed use of targeted sequencing to group *E. coli* clonotypes  
87 from patient's direct urine samples (Tchesnokova, Billig et al. 2013), or to detect ciprofloxacin  
88 resistance markers (Stefan, Koehler et al. 2016), resulting in antimicrobial susceptibility data  
89 and improvement in clinical outcome prediction. Finally, agnostic (unbiased, metagenomic)  
90 sequencing shows promise as a diagnostic of last resort where no other diagnostic can  
91 determine the infectious microorganism, such as the successful ID-NGS diagnosis of leptospira  
92 infection with resulting positive outcome for the patient (Wilson, Naccache et al. 2014).

93 ID-NGS is finding application across the infectious disease space; however, several  
94 studies document the continued need for NGS research and database curation to facilitate  
95 adoption in the clinical setting (Schlaberg, Chiu et al. 2017). Perhaps the best example,

96 Afshinnekoo et. al. showed ID-NGS misidentification of anthrax and plague in the NYC subway  
97 system based on low quality reference genomes (Afshinnekoo, Meydan et al. 2015). A follow-up  
98 erratum by the same group (Afshinnekoo, Meydan et al. 2015) revealed the lack of evidence for  
99 biothreat organisms in these samples. This erratum attributed the anthrax misidentification to  
100 poor reference genomes leading to misattribution of toxin genes when using metagenomic data  
101 analysis tools. This lack of proper reference genomes is pervasive and represents significant  
102 knowledge gaps in public resources, thus emphasizing the necessity for targeted development  
103 of representative, accurate and well curated microbial reference genome sequences. Additional  
104 studies showed that effective use of agnostic sequencing technology, either for infectious  
105 disease identification or exclusion of infectious etiologies, is directly related to the availability of  
106 quality controlled whole-genome reference sequences (Greninger, Messacar et al. 2015,  
107 Naccache, Peggs et al. 2015, Somasekar, Lee et al. 2017). Significant efforts are still required for  
108 ID-NGS technology to transition into a routine clinical diagnostic. To facilitate this transition,  
109 prominent groups and researchers in the field have outlined steps required for proper ID-NGS  
110 use in the clinic (Gargis, Kalman et al. 2016, Schlaberg, Chiu et al. 2017, Simner, Miller et al.  
111 2017).

112 In 2016, the FDA published a draft guidance for ID-NGS devices soliciting feedback on a  
113 potential regulatory pathway for targeted and pathogen-agnostic NGS diagnostic applications.  
114 This draft guidance proposed a novel regulatory strategy for ID-NGS device validation allowing  
115 wet-lab validation of an assay-specific subset of clinical samples to determine the assay

116 preliminary diagnostic performance in combination with *in silico* validation of additional  
117 sequence targets. This *in silico* validation would entail the use of raw sequence data as an input  
118 into bioinformatic algorithms that allow a head-to-head comparison to reference genomes  
119 from the FDA dAtabase for Regulatory-Grade micrObial Sequences (FDA-ARGOS). Figure 1A  
120 illustrates this proposed novel composite reference method (C-RM). By comparison, the current  
121 regulatory paradigm relies on comparing performance of a new test device to FDA Benchmarks  
122 that are either reference standards or non-reference standards (predicate or comparison  
123 method) (See <https://www.fda.gov/RegulatoryInformation/Guidances/ucm071148.htm>).

124 To reduce some of the data generation load in clinical trials, we established and  
125 populated the FDA-ARGOS database with quality controlled microbial sequences as a tool for *in*  
126 *silico* target sequence validation. In this context, *in silico* target sequence validation is part of  
127 the C-RM method focused on evaluating dry lab components (bioinformatic analysis pipelines  
128 and databases) of ID NGS diagnostic assays. Using raw sequence data from the ID-NGS test  
129 device, *in silico* comparison of results obtained with the assay in-house database to results  
130 when using FDA-ARGOS will evaluate device bioinformatic analysis pipelines and report  
131 generation while eliminating the need for additional sample testing with a gold standard  
132 comparator (current FDA benchmarks). Overall, we anticipate the use of the C-RM method  
133 based on assay-specific subsets of clinical samples and/or microbial reference materials  
134 (MRMs) for wet lab validation and FDA-ARGOS *in silico* target sequence validation to generate  
135 scientifically valid evidence for understanding the performance of ID NGS diagnostic assays.

136

137           This manuscript provides our rationale and quality metrics for the FDA-ARGOS genome  
138 database initiative, outlines the need for genome gap filling in the public domain and proposes  
139 the utility of the FDA-ARGOS database resource as a novel *in silico* validation strategy for ID-  
140 NGS diagnostics.

141

142

### 143 **Materials and Methods:**

144

145 **FDA-ARGOS database genome deposition:** Using previously identified microbe(s), nucleic acid  
146 was extracted for library preparation and sequencing. Next, microbial nucleic acids are  
147 sequenced, and de novo assembled using Illumina and Pac Bio sequencing platforms at the  
148 Institute for Genome Sciences at the University of Maryland (UMD-IGS). The assembled  
149 genomes were quality controlled by an ID-NGS subject matter expert working group consisting  
150 of FDA personnel and collaborators with all passing data deposited in NCBI databases. Follow  
151 this link (<https://www.fda.gov/argos>) for full background, collaborators and FDA-ARGOS  
152 genome status. Supplemental Table 1 lists all FDA-ARGOS genomes with accessions and  
153 statistics used in this manuscript.

154

155 **Bacterial reference genome sequencing and assembly:** A hybrid sequencing approach (Koren,  
156 Schatz et al. 2012) based on long and short read NGS technology was selected using Illumina  
157 and PacBio NGS technologies to generate high quality bacterial genome sequences. Sufficient  
158 and high molecular weight genomic starting material was needed for both technologies. Sets of  
159 bacterial libraries were multiplexed on the Illumina PE HiSeq4000 using the 150bp paired-end  
160 run protocol with 24 – 48 isolates per lane. The coverage threshold was set at 300x to ensure  
161 sufficient read depth was achieved from short read NGS technology for high quality assembly  
162 generation. In addition, sets of bacterial libraries were run on the PacBio RS II P6-C4 with at  
163 least 1 SMRT cell per bacterial genome. The coverage threshold was set at 100x to ensure  
164 sufficient and economically feasible read depth was achieved from long read NGS technology  
165 for high quality assembly generation. The data were assembled both separately and in  
166 combination using a series of assembly tools, including SPAdes(Bankevich, Nurk et al. 2012),  
167 Canu (Koren, Walenz et al. 2017), HGAP (Chin, Alexander et al. 2013) and Celera Assembler  
168 (Berlin, Koren et al. 2015). Pilon (Walker, Abeel et al. 2014) was used for polishing of data.  
169 Manual curation was performed to achieve optimal assembly and consensus calling.

170

171 **Viral reference genome sequencing and assembly:** Viral genome sequencing included shotgun,  
172 amplicon, and 5'/3' RACE sequencing methods to generate full-length viral genome sequences.  
173 Sufficient and high quality genomic starting material was needed for all three approaches.  
174 Amplicon sequencing with 48 – 96 overlapping amplicons was used to generate deep coverage



175 of known regions of the genome and was used to evaluate quasi-species in each isolate. Rapid  
176 amplification of cDNA Ends (RACE) was used to finish the 5' and 3' ends, and a shotgun  
177 approach generated data from all RNAs present in the sample without the level of bias present  
178 in the amplicon approach. Sets of viral libraries from all three approaches were multiplexed on  
179 the Illumina MiSeq using the 300bp paired-end run protocol. The coverage threshold was set at  
180 100x to ensure two times amplicon coverage across the genome. The shotgun, amplicon and  
181 RACE data were assembled both separately and in combination using a series of assembly tools,  
182 including SPAdes (Bankevich, Nurk et al. 2012) and Celera Assembler (Berlin, Koren et al. 2015).  
183 Manual curation was performed to achieve optimal assembly and consensus calling.

184

185 **Calculation of FDA-ARGOS genome assembly quality control statistics:** Coverage statistics  
186 were calculated for each of the FDA-ARGOS genome assemblies. Illumina coverage and PacBio  
187 coverage were calculated separately. Illumina short reads were first aligned to the assembly  
188 consensus sequence using Bowtie2 (Langdon 2015). Illumina coverage was then calculated  
189 using samtools (Li, Handsaker et al. 2009) on the resulting sam file. PacBio reads were aligned  
190 to the assembly consensus sequence using BLASR (Chaisson and Tesler 2012). PacBio coverage  
191 was then calculated using samtools (Li, Handsaker et al. 2009) on the resulting sam file. Total  
192 coverage was calculated by adding the PacBio coverage and Illumina coverage at every base  
193 pair location in the assembly consensus sequence.

194

195 **FDA-ARGOS genome annotations:** Genomes were annotated with NCBI's annotation tools to  
196 streamline the process (Angiuoli, Gussman et al. 2008, Brister, Bao et al. 2010, Klimke,  
197 O'Donovan et al. 2011, Tatusova, DiCuccio et al. 2016, Hatcher, Zhdanov et al. 2017). Bacterial  
198 sequences were annotated with NCBI's Prokaryotic Genome Annotation Pipeline (PGAP) that  
199 combines ab initio gene prediction algorithms with homology based methods. Viral sequences  
200 were aligned with their most similar NCBI RefSeqs (NC\_002549, NC\_014372, NC\_006432,  
201 NC\_014373, NC\_004162, NC\_004161, NC\_003899, NC\_001449, NC\_001544, NC\_035889), using  
202 the Geneious alignment tool in the Geneious platform (Kearse, Moir et al. 2012). The setting to  
203 automatically determine detection was used, and the other parameters were set to the  
204 defaults. Gene, CDS, and mature peptide annotations from the RefSeqs were transferred to the  
205 sequences, beginning and end positions were verified for homology, and the sequences were  
206 manually reviewed for unexpected stop codons or regions of high dissimilarity. The RefSeqs  
207 used have had their annotation reviewed by NCBI curators based on available literature, and in  
208 several cases, the annotations were performed in collaboration with researchers familiar with  
209 the viruses.

210

211 **Clinical sample collection and preparation:** Clinical and mock-clinical sample testing was  
212 conducted to demonstrate the utility of FDA-ARGOS. Fifteen de-identified human serum  
213 samples that were Ebola virus (EBOV) Makona positive were received from Sierra Leone; these  
214 samples were determined by the USAMRIID Office of Human Use and Ethics to be Not Human

215 Subject Research (HP-09-32). All samples were collected and de-identified in Sierra Leone at the  
216 Kenema Government Hospital, and the samples had indirect identifiers upon receipt. Presence  
217 of virus for the human samples was determined using the previously established real-time RT-  
218 PCR assay (Trombley, Wachter et al. 2010). Samples were run in duplicate using 5 $\mu$ l of purified  
219 RNA on the LightCycler 480 (Roche Diagnostics Corporation). A positive sample was defined as  
220 having a quantitation cycle (Cq) value of <40 cycles with duplicate positive real-time PCR results  
221 (Table 1B).

222  
223 Ten de-identified human serum samples that were suspected Bundibugyo virus positive were  
224 received from the Democratic Republic of Congo (DRC). These samples were determined by the  
225 USAMRIID Office of Human Use and Ethics to be Not Human Subject Research (HP-12-15).  
226 Presence of virus for the human samples was determined using the previously established  
227 Bundibugyo virus real-time RT-PCR assay (Trombley, Wachter et al. 2010). Samples were run in  
228 duplicate using 5 $\mu$ l of purified RNA on the LightCycler 480 (Roche Diagnostics Corporation). A  
229 positive sample was defined as having a quantitation cycle (Cq) value of <40 cycles (Table 1B).

230  
231 One clinical *Enterococcus avium* from Children's Hospital was used for this study and  
232 maintained at USAMRIID through the Unified Culture Collection (UCC) system. Following  
233 overnight growth of *E. avium*, (~16 hrs), a single, isolated colony was chosen and inoculated  
234 into tryptic soy broth (ThermoFisher, Waltham MA). A glycerol stock was made from the

235 overnight culture and colony counts were performed concurrently to determine the CFU/mL of  
236 the stock organism.

237

238 **Metagenomic and isolate shotgun sequencing:** The *Enterococcus avium* sample

239 SAMN04327393 was cultured on blood agar plates or in tryptic soy broth (ThermoFisher,  
240 Waltham MA). Samples were spiked to a final concentration of  $10^5$  CFU/ml in water or whole  
241 blood matrix (BioreclamationIVT, Baltimore, MD) and 100 $\mu$ l was extracted using the Qiagen EZ1  
242 viral kit (Qiagen, Valencia, CA) according to the manufacturer's instructions. DNA concentration  
243 was quantified utilizing Qubit dsDNA BR assay kit (ThermoFisher). DNA samples were prepared  
244 for sequencing on the MiSeq platform utilizing the Nextera XT DNA library preparation kit  
245 according to the manufacturer's instructions (Illumina, San Diego, CA). Library preparations  
246 were quantified and normalized utilizing the KAPA library quantification kit (Kapa Biosystems,  
247 Wilmington, MA) and sequenced on the MiSeq platform using the 2x150 cycle sequencing kit  
248 (Illumina). Sequencing reads were analyzed using CLC Genomic Workbench (CLC Bio,  
249 Cambridge, MA). For metagenomic analysis, paired end reads were trimmed utilizing a quality  
250 trim of 0.05 and reads below 50bp in length were removed from further analysis. Trimmed  
251 reads were then mapped to *E. avium* assembly GCF\_000407245.1 and *H. sapiens* assembly  
252 GCA\_000001405.27. Mapping parameters were as follows: mismatch costs=2, insertions  
253 costs=3, deletion costs=3, length and similarity fraction = 0.8.

254

255 **Targeted molecular inversion probe sequencing (MIPS):** The Bundibugyo virus (BDBV) and  
256 Ebola virus (EBOV) Makona clinical data samples were run using the previously described MIPS  
257 approach (Koehler, Hall et al. 2014) to capture a targeted sequence into a circular  
258 oligonucleotide. A PCR reaction and subsequent NGS on the Illumina MiSeq (2x150) amplified  
259 and identified the captured sequence using CLC genomics workbench (CLC Bio, Cambridge, MA)  
260 read mapping back to the reference genome (EBOV (GenBank # NC\_002549), BDBV (GenBank #  
261 NC\_014373)). The percent reads classified as Bundibugyo virus or EBOV Makona was reported.  
262 The threshold for positive calls was determined by the no template control (NTC). For the MIPS  
263 approach, the remaining reads are non-specific or “junk”.

264

265 **Mock Clinical Diagnostic Evaluation:** The MIPS assay was evaluated for diagnostic performance  
266 across 148 blinded samples. The limit of detection (LOD) was determined through a preliminary  
267 titration of EBOV Zaire in TRIzol starting at  $10^8$  plaque forming units (pfu)/ml down to  $10^2$   
268 pfus/ml and then run in triplicate. The concentration where all three replicates yielded positive  
269 results was confirmed as the LOD across 40 replicates at that concentration. EBOV (Kikwit  
270 R4317a) in TRIzol LS was diluted to 10X ( $1.0E+06$  pfu/ml), 5X ( $5.0E+05$  pfu/ml) and 1X ( $1.0E+05$   
271 pfu/ml) LOD in triplicate in matrix also containing TRIzol LS. Nucleic acid was extracted using  
272 400 $\mu$ l of each sample, along with 14 negative serum samples, on the EZ1 Virus 2.0 kit and  
273 eluted in 60 $\mu$ l. Presence of virus was determined with an established real-time PCR assay in  
274 triplicate for each extracted sample. Extracted RNA was amplified from 5  $\mu$ l total nucleic acid

275 using the Quantitect Whole Transcriptome Amplification Kit (Qiagen) and quantified with the  
276 Qubit dsDNA Broad Range Assay Kit. A total of 50ng cDNA was added into the MIP protocol.  
277 Library preparation was performed on the Apollo instrument using the PrepX Complete ILMN  
278 32i DNA kit and Illumina TruSeq dual Indices. All samples were sequenced on the Illumina  
279 MiSeq using the 300 cycle kit. Sixteen samples were spiked at 10X, 5X and 1X LOD. For the mock  
280 clinical evaluation, 48 positive and 100 negative (matrix only) samples were run as described  
281 above. Threshold cutoffs for positive samples were 2X signal to noise ratio (SNR). All diagnostic  
282 performance statistics were calculated on [https://www.medcalc.org/calc/diagnostic\\_test.php](https://www.medcalc.org/calc/diagnostic_test.php).

283

284 **Short Read Classification Using MegaBLAST Tool:** The quality of the short reads was checked  
285 with FastQC. No quality trimming was conducted. We selected 100,000 short reads randomly  
286 from each of the samples (140,000 for mock clinical). The MegaBLAST function of blast+ 2.7.1  
287 installed on FDA HPC infrastructure (<https://www.ncbi.nlm.nih.gov/books/NBK153387/>) was  
288 used to taxonomically classify the short reads using the default parameters and three  
289 databases: Algorithm Standard Database (NCBI Nt), Algorithm Standard Database and FDA-  
290 ARGOS and FDA-ARGOS alone. NCBI Nt was downloaded and constructed on 9/25/2017. The  
291 FDA-ARGOS database was constructed with FDA-ARGOS genomes (Supplemental Table 1,  
292 SAMN04327393 was excluded from the database because this reference genome was  
293 developed from the same isolate that was used as spike in material for use case 1) using the  
294 makeblastdb command. The Algorithm Standard Database and FDA-ARGOS database was

295 constructed by aggregating the NCBI Nt database and FDA-ARGOS database. Default options  
296 were used to build the databases. For this study, the taxon associated with the first reported  
297 alignment was used as the taxonomic label for each read. Original MegaBLAST results were  
298 summarized to report the number of reads associated with each unique NCBI taxonomy ID  
299 called.

300

301 **Short Read Classification Using Kraken Tool:** The quality of short reads was checked with  
302 FastQC. No quality trimming was conducted. We subsampled 300,000 short reads uniformly  
303 from each of the samples. Kraken 1.0 (Wood and Salzberg 2014) , installed on FDA HPC  
304 infrastructure, was used to assign a taxonomic label to each short read using default  
305 parameters and three databases: Algorithm Standard Database (NCBI Nt), Algorithm Standard  
306 Database and FDA-ARGOS and FDA-ARGOS alone. NCBI Nt was downloaded and constructed on  
307 10/5/2017. The FDA-ARGOS database was constructed with FDA-ARGOS genomes  
308 (Supplemental Table 1, SAMN04327393 was excluded from the database because this  
309 reference genome was developed from the same isolate that was used as spike in material for  
310 use case 1) using Kraken-build command. The Algorithm Standard Database and FDA-ARGOS  
311 database was constructed with both the NCBI Nt database and the FDA-ARGOS genomes.  
312 Default options were used to build the databases. For this study, the taxon associated with the  
313 first reported alignment was used as the taxonomic label for each read. Original Kraken results

314 were summarized to report the number of reads associated with each unique NCBI taxonomy  
315 ID called.

316

317 **Short read Classification Using LMAT:** The quality of the short reads was checked with FastQC.  
318 No quality trimming was conducted. LMAT version 1.2.6 (available for download at  
319 [sourceforge.net/lmat](https://sourceforge.net/lmat), (Ames et al., 2015)), installed on Lawrence Livermore National  
320 Laboratory (LLNL) HPC infrastructure was used to assign a taxonomic label to each short read  
321 with a minimum score setting of 0.5. Match scores are calculated per read, by fitting a random  
322 null model created by simulating 1 GB of random sequence for each model dependent on read  
323 length and GC content. Three databases, the Algorithm Standard Database (LMAT DB), the  
324 stand-alone FDA-ARGOS database (Supplemental Table 1, SAMN04327393 was excluded from  
325 the database because this reference genome was developed from the same isolate that was  
326 used as spike in material for use case 1) and an aggregated database consisting of both the  
327 LMAT DB database and the stand-alone FDA-ARGOS database were used. LMAT results were  
328 summarized to report the number of reads associated with each unique NCBI taxonomy ID.

329

330

331 **Results:**

332

333 *Filling gaps in public resources with targeted reference genomes*



334

335           In 2013, FDA in collaboration with the Department of Defense (DoD) and the National  
336 Center for Biotechnology Information (NCBI) assessed the quality and diversity of sequenced  
337 microbial genomes present in public databases. A majority of pathogens appeared to be  
338 represented by multiple entries, however, many of these genomes were incomplete or of  
339 unknown quality. In fact, a thorough examination of the entire public domain revealed some  
340 pathogens were underrepresented or completely absent. Our 2013 review, supported by  
341 several publications (Schatz and Langmead 2013, Land, Hyatt et al. 2014, Land, Hauser et al.  
342 2015), revealed biased phylogenetic coverage usually attributable to research funding for  
343 specific microbial model organisms. At the time, NCBI GenBank covered less than 8,000  
344 bacterial and archaeal genome sequences with at least half submitted by the four largest  
345 genome sequencing centers: Broad Institute, DOE Joint Genome Institute, Institute for Genome  
346 Sciences and TIGR/JCVI. Additionally, many sequences lacked accompanying metadata and raw  
347 read information. These issues provided the impetus for *de novo* generation of FDA-sponsored  
348 reference sequences of the highest quality achievable using state-of-the-art genomic  
349 sequencing technologies (Koren, Schatz et al. 2012). With this effort, FDA intended to establish  
350 quality control metrics for microbial genomes that could be used in ID-NGS test validation. Only  
351 genomes with the highest technically achievable quality would qualify as regulatory-grade  
352 genomes. Factors essential to reach that goal were: 1) knowledge of the technology used to  
353 generate the sequences, 2) access to raw sequence information to reproduce the data, and, 3)

354 access to relevant metadata. Perhaps the most significant missing piece of information for  
355 previously generated reference genomes was the lack of an independent reference method  
356 that reliably linked the microbial organism identification to the sequence data. In this context,  
357 qualification of microbial reference genomes requires organism identification with a  
358 recognized reference method as this remains a primary requirement for validation of a new  
359 diagnostic device.

360 FDA, DOD, NCBI and other agencies using scientific literature, a phylogenetic data  
361 mining approach, and FDA microbial species-specific guidance documents identified more than  
362 1000 gaps in public microbial genomic repositories. We prioritized these gaps and selected  
363 biothreat microorganisms, common clinical pathogens and closely related species (See  
364 Supplemental Materials for the organism gap list). The primary objective of this regulatory  
365 science research and tool development effort centered on the generation of an initial set of  
366 2000 quality-controlled microbial FDA-ARGOS reference genomes . These genomes are  
367 generated with a hybrid assembly approach using short and long read sequencing technologies  
368 (Koren, Schatz et al. 2012). An initial collection criterion focused on sequencing at least 5  
369 diverse isolates per species to cover temporal and spatial genome plasticity and initiate the  
370 construction of a regulatory-grade microbial genome model.

371

372 *FDA-ARGOS, what's that?*

373

374 FDA and collaborators established the publicly available database, FDA dAtabase for  
375 Regulatory-Grade micrObial Sequences (FDA-ARGOS), to fill these defined gaps for genomic  
376 sequences. Here, we present the first subset of 487 FDA-ARGOS genomes with NCBI accessions  
377 (Figure 2, Supplemental Table 1). Of the 487 isolates, 88.3 percent were bacteria, 11.1 percent  
378 were viruses and 0.6% were eukaryotes, representing 189 different taxa. In total, 81.9 percent  
379 of genomes were of clinical origin with the remaining 18.1 percent environmental genomes  
380 from closely related species near-neighbors (Supplemental Table 2). Over 500 isolates are  
381 currently being sequenced and at different stages in the FDA-ARGOS genome generation  
382 pipeline.

383 Use of advanced sequencing technologies (Koren, Schatz et al. 2012) helped define the  
384 characteristics for regulatory-grade genomes. Specifically, Figure 1B provides a summary of  
385 required FDA-ARGOS metrics to support a determination of regulatory-grade genome. All FDA-  
386 ARGOS genomic submissions demonstrated: 1) organism identification prior to sequencing by a  
387 recognized independent reference method, 2) sequence generation with at least two  
388 sequencing methodologies (e.g., long read and short read NGS), and, 3) *de novo* assembly with  
389 high-depth of base coverage. Each microbial isolate assembled genome sequence conformed to  
390 a minimum of 95 percent coverage with 20X depth at every position while also providing  
391 concordant NCBI taxonomy-specific average nucleotide identity (ANI) thresholds for microbial  
392 organism identification (Ciufo, Kannan et al. 2018) with independent identification methods. All  
393 FDA-ARGOS samples were concordant between *de novo* sequencing identification and

394 independent organism identification method (Supplemental Table 2 lists independent  
395 identification method data).

396 As mentioned above, hybrid error-correction with long and short read sequencing  
397 technology was considered for establishing minimum FDA-ARGOS regulatory grade data  
398 requirements. Figure 1C outlined these criteria and included sample name, 10 meta data fields  
399 (based on NCBI BioSample submission requirements), raw reads, assemblies with coverage,  
400 N50, L50 and annotations. Importantly, FDA-ARGOS genomes are tied to a minimum of 10  
401 critical sample metadata fields (Figure 1D): independent organism confirmation by recognized  
402 reference method, culture collection, and, the following required NCBI BioSample fields:  
403 organism, strain, isolation source, host, collected by, taxonomy ID, contact and package  
404 information. Supplemental Table 2 shows metadata coverage metrics for all 487 FDA-ARGOS  
405 genomes. The 10 sample metadata fields are 100 percent completed and available throughout  
406 the sample set with 5 additional metadata metrics are recommended, such as geographic  
407 location, collection date, host disease, host sex and host age (BioSample documentation  
408 <https://www.ncbi.nlm.nih.gov/biosample/docs/attributes/>). In terms of clinical representation,  
409 81.9 percent of clinical samples in the collection are associated with known phenotype/host  
410 disease.

411 Critical for the designation of genomes as ‘regulatory-grade genomes’, was the  
412 institution of quality control metrics for all aspects of the genome generation. To objectively  
413 identify such quality control metrics, we performed internal quality control assessments of all

414 487 genome assemblies (See methods for calculation of FDA-ARGOS genome assembly quality  
415 control statistics, Supplemental Table 1). Figure 3 shows the quality of FDA-ARGOS genome  
416 assemblies compared to the representative 2013 NCBI GenBank database and the  
417 representative 2018 NCBI GenBank database. Both, the 2013 and 2018 NCBI database captures  
418 held up to 50 NCBI assemblies for each species within the FDA-ARGOS database from the  
419 respective year. In relative number of assemblies, 2018 NCBI database contained 3535 while  
420 the 2013 contained 1617. Overall, we observed higher quality in the FDA-ARGOS genome  
421 dataset for the coverage, N50, and L50 quality assembly metrics compared to the 2013 and  
422 2018 NCBI GenBank public genome dataset (Figure 3 A, B and C respectively). Figure 3D  
423 demonstrated that only 675 out of the 3535 2018 NCBI GenBank assembled genomes, or 20  
424 percent, showed comparative assembly quality to FDA-ARGOS genome sequences when  
425 considering one of the reported assembly quality metrics. More importantly, when considering  
426 all quality control assembly metrics, only 11 out of the 3535 2018 NCBI GenBank assembled  
427 genomes, or 0.3 percent, showed comparable quality to FDA-ARGOS genome assemblies.

428 We expect refinement of the quality metrics for ‘regulatory-grade’ genome status  
429 (Figure 1B) as we continue to populate the FDA-ARGOS with additional quality-controlled  
430 genomes; therefore, we established the requisite that all genomes should be available publicly.  
431 Deposition for all FDA-ARGOS genomes requires that raw reads, assembled genomes and  
432 associative metadata are publicly available (<https://www.ncbi.nlm.nih.gov/bioproject/231221>).  
433 (Check <https://fda.gov/argos> for additional background information and updated genomes).

434

435 *FDA-ARGOS fills critical gaps in public sequence repositories - Use Case 1: Enterococcus avium*

436

437           Several regulatory science research considerations arose during the process of  
438 generating FDA-ARGOS genomes, including the initial impetus for this effort, gap filling. Our  
439 first use case documented the importance of genome gap filling with FDA-ARGOS quality-  
440 controlled genomes, and the impact of lack of publicly available genomes for medically  
441 important microbes on potential diagnostic applications. Specifically, we tested whether the  
442 addition of quality-controlled reference sequences into the public repositories impacted the  
443 NGS pathogen detection of a metagenomic shotgun sequencing approach of a mock clinical *E.*  
444 *avium*-spiked human blood sample at clinically relevant titers. An isolate from reference  
445 genome SAMN04327393, which was removed from reference databases for data analysis, was  
446 used as a mock clinical *E. avium* sample. Initial read mapping using CLC Genomics and *E. avium*  
447 sequences from publicly available databases as a reference demonstrated *de novo* assembly of  
448 *E. avium* data was not possible due to only an average 424.4 mapped paired end reads  
449 (Supplemental Table 7). For frame-of-reference, we would need over 30,000 reads to *de novo*  
450 assemble an entire genome of approximately 5 Mb at 1X coverage, assuming a read size of 150  
451 bp and perfect quality of each generated read at all positions.

452           Subsequent bioinformatics data analysis of the *E. avium* metagenomics shotgun paired-  
453 end reads data showed the critical gap filling and utility of the FDA-ARGOS database resource.

454 We analyzed the effect of genome gap filling with MegaBLAST (Morgulis, Coulouris et al. 2008)  
455 and Kraken (Wood and Salzberg 2014) by determining the number of *E. avium* reads classified  
456 from the mock clinical human blood sample with and without FDA-ARGOS genomes used in the  
457 respective bioinformatics tools reference databases. Intuitively, a majority, over 98 percent, of  
458 approximately 12 million paired-end reads for each replicate sample mapped against the  
459 human genome with only 2 percent or less mapping to non-human sequences with both  
460 algorithms (Figure 4A). In contrast, application of MegaBLAST and Kraken with FDA-ARGOS  
461 alone yielded zero human reads due to the lack of human reference in that database. Reads  
462 classified as *E. avium* ranged from an average 3829 and 840 when FDA-ARGOS genomes were  
463 added to the algorithm reference database compared to an average 29 and 0 reads when these  
464 genomes were absent for MegaBLAST and Kraken, respectively (Figure 4B, Supplemental Table  
465 4). Interestingly, while *E. avium* genomes were available in the NCBI Nt database and part of  
466 the read classification for MegaBLAST analyses, positive ID-NGS identification required the  
467 addition of quality-controlled FDA-ARGOS reference genomes. MegaBLAST tool with FDA-  
468 ARGOS data as the standalone reference database generated the largest effect with an  
469 additional 1495 *E. avium* reads classified. MegaBLAST classified additional reads with the stand-  
470 alone FDA-ARGOS database most likely because the quality-controlled *E. avium* genomes were  
471 not mixed with lower quality genomes in the standard algorithm database as competition with  
472 lower quality and closely related genomes removed.

473           Finally, we performed *E. avium* isolate shotgun sequencing without clinical matrix to  
474 obtain sufficient data to illustrate the critical nature of having quality-controlled reference  
475 genomes. Using the aforementioned bioinformatics tools, data analysis showed the impact of  
476 read classification solely focused on *E. avium* and determined if the addition of FDA-ARGOS  
477 genomes to public databases affected read mapping. Figure 4C shows that addition of FDA-  
478 ARGOS *E. avium* reference genomes significantly increased read classification performance  
479 based on the number and percent of *E. avium* reads classified (Figure 4C, Supplemental Table  
480 5). On average, for *E. avium* isolate shotgun sequencing, 8,406,630 reads out of a total 12  
481 million reads classified as *E. avium* when the FDA-ARGOS database resource upon addition to  
482 the algorithm standard reference database (NCBI Nt) compared to 25,800 reads without FDA-  
483 ARGOS added. Amalgamation of FDA-ARGOS genomes into standard sequence reference  
484 databases resulted in *E. avium* contributing between 84 to 96 percent of the total reads  
485 classified (Figure 4C). Interestingly, top hits from the MegaBLAST tool using NCBI Nt database  
486 (containing 4 *E. avium* genomes but not at regulatory grade quality, Supplemental Table 3)  
487 showed over 10 percent of total classified reads mapped to 'Bos Taurus' or '*Enterococcus*  
488 *faecium*'. These top hits were potentially database contaminants and illustrate the risk of using  
489 non-curated databases in ID-NGS diagnostics. Data analysis with the Kraken tool and the  
490 algorithm standard reference database (NCBI Nt) resulted in 0 mapped reads because the  
491 Kraken tool reference database lacked *E. avium* genomes (Figure 4C).



492 For future benchmarking efforts of bioinformatics tools, we provide all *E. avium* Data  
493 Sets (Supplemental Material).

494

495 *In Silico Comparison: Regulatory-grade genomes are sufficient for Ebolavirus Target Sequence*

496 *Validation – FDA-ARGOS Use Case 2*

497

498 A major incentive for the development of FDA-ARGOS was to enable and promote  
499 innovation for ID-NGS medical devices. Through the process of populating the FDA-ARGOS  
500 database, the concept of partial *in silico* validation, rather than completely empirical validation  
501 of clinical trial samples with an independent gold standard reference method, matured. We  
502 chose FDA-ARGOS Ebola reference sequences (Supplemental Table 1) and a targeted ID-NGS  
503 assay, the Ebola virus molecular inversion probes (MIPS), to evaluate the application of FDA-  
504 ARGOS as an *in silico* target sequence validation tool. Table 1 showed the diagnostic  
505 performance of the MIPS ID-NGS assay with clinical Bundibugyo virus and EBOV Makona  
506 samples reported as a more sensitive assay, EBOV Real-Time PCR (RT-PCR) assay (Trombley,  
507 Wachter et al. 2010). When assessing 10 clinical Bundibugyo virus and 15 clinical EBOV Makona  
508 samples, concordant real-time PCR and MIPS positive results ranged from 9 out of 10 clinical  
509 samples (Table 1A) to 6 out of 15 (Table 1B), respectively. Intuitively, lower quantitation cycle  
510 ( $C_q$ ) values correlated with higher MIPS read classification, suggesting the capability of ID-NGS  
511 to detect organisms was dependent on the starting concentration of the target genomic

512 material. MIPS false negative calls for low target analytes suggested that complete *in silico*  
513 validation is an unrealistic approach for clinical trials without comparison to some gold standard  
514 reference method, in this case real-time PCR.

515 Consistent concordance between the benchmark RT-PCR assay, the MIPS test device and  
516 the FDA-ARGOS *in silico* target sequence validation was important for establishing confidence in  
517 considering *in silico* comparison method for clinical sample ID calling. To test this assumption,  
518 we used three bioinformatics tools, MegaBLAST (Morgulis, Coulouris et al. 2008), Kraken (Wood  
519 and Salzberg 2014) and LMAT (Ames, Hysom et al. 2013) to evaluate the proposed *in silico*  
520 target sequence validation method (Figure 1A), and to verify the potential for using *in silico*  
521 comparison without any empirical validation. MegaBLAST and Kraken analyses of raw  
522 sequence data for Bundibugyo virus samples using the three different read classification tools  
523 in combination with FDA-ARGOS as the reference genome database showed complete  
524 agreement for MIPS and *in silico* calls (Table 1A). Because the *in silico* comparison missed the  
525 classification call against the gold standard PCR benchmark test for a sample with low analyte  
526 levels (1 false negative result for the *in silico* validation), we performed a more in depth analysis  
527 of the additional EBOV Makona samples across three bioinformatics tools, MegaBLAST, Kraken  
528 and LMAT (Table 1B). These analyses showed similar results to the Bundibugyo virus data at  
529 100% agreement with test device, but only for samples with low  $C_q$  or high input concentrations  
530 of the target organism. Additional analyses comparing results for each bioinformatics tool  
531 reference databases with and without FDA-ARGOS genomes added, produced similar results

532 demonstrating that FDA-ARGOS alone was sufficient for *in silico* comparison (Supplemental  
533 Table 6). Overall, these data suggested *in silico* sequence comparison would be completely  
534 reliant on the inherent sensitivity of the sequencing assay to generate sequence read data for  
535 comparison, therefore Composite Reference Method (C-RM) (combining *in silico* sequence  
536 comparison with a wet lab validation challenge) is necessary for full validation of the test ID-  
537 NGS device. Figure 1A illustrates the proposed novel C-RM, highlighting this need for empiric  
538 assessment of an ID-NGS assay-specific subset of samples or well defined microbial reference  
539 materials.

540 Evaluation of the clinical samples suggested a need for benchmarking ID-NGS assays to  
541 currently implemented reference methods, thus the application of the C-RM. To document the  
542 application of MIPS Ebola Makona ID-NGS assay benchmarking, we performed a mock clinical  
543 trial to assess the assay-specific wet-lab subset evaluation as part of the proposed C-RM.  
544 Initially, we performed a preliminary limit of detection (LOD) evaluation to determine the scope  
545 of the mock clinical evaluation. These experiments showed a preliminary LOD of  $10^5$  with linear  
546 dose response correlation to EBOV input across the titration (Supplemental Table 8). An  
547 additional 40 positive replicates performed on two independent days, two independent runs  
548 confirmed the LOD at  $10^5$  pfu/ml for EBOV. This concentration formed the basis for spike-in  
549 levels of the mock clinical trial. From a total of 148 samples tested, 48 constituted positive  
550 spiked samples with 16 at high (10x LOD), 16 at medium (5x LOD) and 16 at 1x LOD for the MIPS  
551 assay (Table 2). In this mock clinical trial, all spiked samples were positive via real-time PCR

552 (data not shown). Only 9 out of 16 samples at 1X LoD for the MIPS assay were positive with 37  
553 out of 48 samples positive across the entire sample set in this analysis. However, the positive  
554 predictive value (PPV) and negative predictive value (NPV) for the MIPS assay were: 97.4% and  
555 90%, respectively at or above the limit of detection with a prevalence of 32.4%. In addition,  
556 Table 2 lists the positive and negative predictive values for prior probabilities of infection from  
557 0-1. The PPV and NPV metrics are important predictive analytics tools to provide performance  
558 characteristics for how the ID-NGS diagnostic test will perform in a clinical context. These data  
559 provide a rationale for developers using partial *in silico* validation when false negative rate is  
560 low.

561 For future benchmarking efforts of bioinformatics tools, we have provided all Ebola Data  
562 Sets (Supplemental Material).

563

564

565

## 566 **Discussion:**

567

568 To encourage innovation and support the infectious disease community, we provide here  
569 the FDA-ARGOS resource as a tool for ID-NGS assay development, reference database and *in*  
570 *silico* target sequence validation as part of a novel Composite Reference Method (C-RM). This  
571 manuscript describes the database, specifically highlighting: 1) the quality metrics of regulatory-

572 grade genomes for database inclusion, 2) benefits of FDA-ARGOS in filling pathogen genome  
573 knowledge gaps for device output, and 3) describes some use cases for FDA-ARGOS.

574 A critical aspect for assessing performance of any diagnostic is the availability of minimum  
575 quality control metrics for data, genomic or otherwise, for validation. Defined here are the FDA-  
576 ARGOS 'regulatory-grade' genome criteria that provide ID-NGS diagnostic assay developers and  
577 the scientific community with traceable and quality-controlled genomes. These high-quality  
578 genomes coupled with a streamlined approach for comprehensive expansion of FDA-ARGOS  
579 beyond the initial 2000 genomes is essential for continued ID-NGS diagnostic assay  
580 development.

581 FDA-ARGOS genome sequencing and research resulted in six broad quality metrics (Figure  
582 1B) defining 'regulatory-grade' genome criteria required for current and future FDA-ARGOS  
583 contributors. All extant genomes in the FDA-ARGOS database (Supplemental Table 1) adhere to  
584 the quality metrics of 95% coverage with 20X depth at every position across the entire  
585 assembled genome. This metric applies to the initial deposition of, minimally, 5 genomes of any  
586 genus/ species added to FDA-ARGOS. These 5 genomes define the "FDA-ARGOS core genome".  
587 After 5 or more regulatory grade genomes per genus/species are available in the database, we  
588 will consider lower threshold metrics for FDA-ARGOS inclusion to capture novel and/or unique  
589 genomes that may be diagnostically informative. Future efforts will apply these metrics to  
590 existing genomic information in the public domain coupled with deep learning methods and

591 artificial intelligence to inform an external genome qualification tool greatly expanding utility of  
592 the FDA-ARGOS database.

593         Lack of high quality reference genomes challenges the accuracy of ID-NGS identification  
594 for queryable microbial pathogen. The genome gap filling use case with regulatory-grade *E.*  
595 *avium* genomes highlights current challenges with infectious disease NGS technology when  
596 using minimal-, non-curated or absent reference databases. The end result potentially leading  
597 to the lack of a diagnostic call or even misdiagnosis. These data were punctuated by two key  
598 findings: 1) *de novo* assembly of the data was not possible due to the low number of reads in  
599 clinical matrix and 2) limited *E. avium* species reference genomes in publicly available databases  
600 made the sample identification almost impossible (Supplemental Table 3). The latter point is  
601 extremely relevant for the intent of ID-NGS for diagnostic applications. In the case presented  
602 here, the top microbial sequence hit did not equate to the microbe of interest due to lack of  
603 representation in the reference database. Intuitively, addition of FDA-ARGOS and relevant  
604 genomes mitigated this issue. In addition, *E. avium* isolate sequencing results showed the  
605 dependency of both classification method (such as MegaBlast and Kraken) and database used.  
606 This last aspect of the *E. avium* use case informed the consideration of the C-RM and opened  
607 the possibility for utilizing a suite of validated bioinformatics tools for *in silico* target sequence  
608 validation.

609         There are two basic contrasting philosophies in circulation regarding genomic information  
610 and ID-NGS: 1) all information, whatever the quality, is useful towards making a diagnosis, the

611 more data the better, with the assumption of diagnosis relying on error correction through  
612 iteration, or, 2) quality-controlled, highly curated genomes are required as a solid foundation,  
613 more information is better, however, diagnostics require quality-controlled genomes to inform  
614 the basis of diagnosis. Experiments and data presented here support the latter of these two  
615 arguments. Specifically, while *E. avium* reference genomes were available in NCBI Nt database  
616 and were part of the read classification for MegaBLAST analyses, positive ID-NGS identification  
617 of *E. avium* required the addition of quality-controlled FDA-ARGOS reference genomes. In  
618 addition, read mapping of isolate shotgun data, without any clinical matrix, showed  
619 indeterminate results for *E. avium* without FDA-ARGOS in contrast to 80 percent of total reads  
620 mapped as *E. avium* upon addition of these regulatory-grade genomes to the reference  
621 database. A similar increase in performance in *E. avium* reads classified resulted when using  
622 FDA-ARGOS *E. avium* reference genomes for metagenomics shotgun data, in whole blood, even  
623 with human reads occupying >98% of sequencing real estate.

624       Quality and coverage of targeted organisms are critical aspects for ID-NGS transition  
625 into the clinical space; however, to foster the transition, new methods are required to lessen  
626 the burden for validating ID-NGS against all queryable pathogens. This manuscript documents  
627 methods for use of FDA-ARGOS reference genomes in *in silico* sequence comparison as part of  
628 the proposed novel C-RM. We showed here that the *in silico* validation of Bundibugyo virus and  
629 Zaire ebolavirus can use FDA-ARGOS genomes as the comparator. For MIPS positive samples,  
630 there was 100 percent concordance between the gold standard real-time PCR comparator, and

631 the *in silico* comparison. This supports the feasibility of implementing this strategy to shorten  
632 future clinical NGS-based assay evaluation studies. A potential mitigation for this issue, where  
633 real-time PCR was more sensitive than the MIPS NGS assay especially at high  $C_q$  values, is the  
634 application of additional enrichment strategies to bring ID-NGS to similar sensitivities as the  
635 gold standard (Briese, Kapoor et al. 2015, O'Flaherty, Li et al. 2018). However, in the current  
636 form, observed lower sensitivity of the MIPS assay compared to real-time PCR shows the  
637 necessity for a C-RM and incorporating additional empirical studies, i.e., an assay-specific  
638 subset of clinical samples going through wet-lab comparison as part of the clinical validation.  
639 Discordant results at high  $C_q$  values highlight the perils of solely applying *in silico* sequence  
640 comparison. Without any empirical evaluation, *in silico* comparison would only provide results  
641 within the sensitivity ranges of the test ID-NGS device without providing the needed benchmark  
642 for sensitivity compared to a gold-standard such as real-time PCR. Therefore, as part of the C-  
643 RM, we demonstrate a preliminary performance assessment against a gold-standard for a  
644 subset of the clinical trial samples with the intent that the remainder of the clinical trial samples  
645 could be validated via *in silico* sequence comparison. Different sample read depths may be  
646 required to achieve the desired identification performance for various organisms. Assay  
647 developers might be required to use an external comparator only for *in silico* validation results  
648 where the test device and *in silico* comparison yielded a discordant result. We envision this C-  
649 RM to be a primary utility of the FDA-ARGOS genome database tool for medical device



650 development. We hope that FDA-ARGOS will spur innovation and expedite regulatory science,  
651 and ultimately enable ID-NGS as a diagnostic to enter the clinic.

652         The FDA-ARGOS reference genome resource is a constantly evolving public database  
653 instance and intended to mature over time with community support and genomic technology  
654 advancements. Continued population and expansion of the FDA-ARGOS database resource will  
655 be required to cover the panoply of infectious microorganisms. In this proposed *in silico*  
656 validation with FDA-ARGOS, the need for comprehensive regulatory-grade genome coverage is  
657 clear, however, no one entity can perform all the needed sequencing. We are therefore  
658 working on a pathway for external genome qualification to streamline and expand FDA-ARGOS  
659 resource as needed. Both the external genome qualification and continued research to apply  
660 this regulatory-grade standard to unculturable and emerging pathogens will be the focus of  
661 future research.

662         Further population and curation of the database will support the success of FDA-ARGOS  
663 and promote adoption by the NGS community. The FDA-ARGOS team openly invites additional  
664 collaborators from the scientific community to assist in filling the gaps in this public resource.  
665 FDA-ARGOS and collaborators are specifically searching for unique, hard to source microbes  
666 such as biothreat organisms, emerging pathogens, and clinically significant bacterial, viral,  
667 fungal, and parasitic genomes. As stated, the goal is to collect sequence information for a  
668 minimum of 5 isolates per species and we solicit any potential collaborators interested in  
669 supplying these 5 isolates for gap-filling to contact and authors of this paper. For more

670 information about contributing samples for UMD-IGS sequencing as part of FDA-ARGOS efforts,  
671 or to qualify existing genomes by the FDA, please email [FDA-ARGOS@fda.hhs.gov](mailto:FDA-ARGOS@fda.hhs.gov).

672

673

#### 674 Acknowledgements

- 675 • This project has been funded with Federal funds from the Office of Counterterrorism  
676 and Emerging Threats, Food and Drug Administration, Department of Health and Human  
677 Services, under Contract No. HHSF223201310109C, HHSF223201510106C,  
678 HHSF223201610073C and the Department of Defense, under Contract No. 224-15-  
679 6506R.
- 680 • This research was supported by the Intramural Research Program of the NIH, National  
681 Library of Medicine.
- 682 • This project was funded by DTRA, Contract No. CB10245.
- 683 • Sample contributions for the initial set of 500 from the U.S. Army Medical Research  
684 Institute of Infectious Diseases, the Department of Defense Critical Reagents Program,  
685 Public Health Agency Canada, Public Health England, the University of Texas Medical  
686 Branch, BC Centre for Disease Control, American Type Culture Collection, Rockefeller  
687 University, FDA-CBER (Maria Rios, Robert Duncan, Rafaele Gusmao), FDA-CFSAN (Eric  
688 Brown, Marc Allard, Maria Hoffman, Cary Pirone), FDA-CVM (Patrick McDermott,  
689 Shaohua Zhao), Children's National Hospital (Joseph Campos, Brittany Goldberg, Chelsie  
690 Geyer), University of Colorado School of Medicine (Thomas Morrison, Sudhakar  
691 Agnihothram)
- 692 • The opinions, interpretations, conclusions, and recommendations contained herein are  
693 those of the authors and are not necessarily endorsed by the U.S. Army.
- 694 • The views expressed here are those of the authors and do not necessarily represent the  
695 views or official position of the FDA, NIH or DOE.

696

697

#### 698 **References:**

- 699 1. Afshinnkoo, E., C. Meydan, S. Chowdhury, D. Jaroudi, C. Boyer, N. Bernstein, J. M. Maritz, D.  
700 Reeves, J. Gandara, S. Chhangawala, S. Ahsanuddin, A. Simmons, T. Nessel, B. Sundares, E. Pereira, E.  
701 Jorgensen, S. O. Kolokotronis, N. Kirchberger, I. Garcia, D. Gandara, S. Dhanraj, T. Nawrin, Y. Saletore, N.

- 702 Alexander, P. Vijay, E. M. Henaff, P. Zumbo, M. Walsh, G. D. O'Mullan, S. Tighe, J. T. Dudley, A. Dunaif, S.  
703 Ennis, E. O'Halloran, T. R. Magalhaes, B. Boone, A. L. Jones, T. R. Muth, K. S. Paolantonio, E. Alter, E. E.  
704 Schadt, J. Garbarino, R. J. Prill, J. M. Carlton, S. Levy and C. E. Mason (2015). "Geospatial Resolution of  
705 Human and Bacterial Diversity with City-Scale Metagenomics." Cell Syst **1**(1): 72-87.
- 706 II. Afshinnakoo, E., C. Meydan, S. Chowdhury, D. Jaroudi, C. Boyer, N. Bernstein, J. M. Maritz, D.  
707 Reeves, J. Gandara, S. Chhangawala, S. Ahsanuddin, A. Simmons, T. Nessel, B. Sundares, E. Pereira, E.  
708 Jorgensen, S. O. Kolokotronis, N. Kirchberger, I. Garcia, D. Gandara, S. Dhanraj, T. Nawrin, Y. Saletore, N.  
709 Alexander, P. Vijay, E. M. Henaff, P. Zumbo, M. Walsh, G. D. O'Mullan, S. Tighe, J. T. Dudley, A. Dunaif, E.  
710 Ennis, E. O'Halloran, T. R. Magalhaes, B. Boone, A. L. Jones, T. R. Muth, K. S. Paolantonio, E. Alter, E. E.  
711 Schadt, J. Garbarino, R. J. Prill, J. M. Carlton, S. Levy and C. E. Mason (2015). "Geospatial Resolution of  
712 Human and Bacterial Diversity with City-Scale Metagenomics." Cell Syst **1**(1): 97-97 e93.
- 713 III. Ames, S. K., D. A. Hysom, S. N. Gardner, G. S. Lloyd, M. B. Gokhale and J. E. Allen (2013).  
714 "Scalable metagenomic taxonomy classification using a reference genome database." Bioinformatics  
715 **29**(18): 2253-2260.
- 716 IV. Angiuoli, S. V., A. Gussman, W. Klimke, G. Cochrane, D. Field, G. Garrity, C. D. Kodira, N.  
717 Kyrpides, R. Madupu, V. Markowitz, T. Tatusova, N. Thomson and O. White (2008). "Toward an online  
718 repository of Standard Operating Procedures (SOPs) for (meta)genomic annotation." OMICS **12**(2): 137-  
719 141.
- 720 V. Arnold, C. (2017). "Source code: Putting metagenomics to the test in the clinic." Nat Med **23**(6):  
721 645-648.
- 722 VI. Bankevich, A., S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, A. S. Kulikov, V. M. Lesin, S. I.  
723 Nikolenko, S. Pham, A. D. Prjibelski, A. V. Pyshkin, A. V. Sirotkin, N. Vyahhi, G. Tesler, M. A. Alekseyev  
724 and P. A. Pevzner (2012). "SPAdes: a new genome assembly algorithm and its applications to single-cell  
725 sequencing." J Comput Biol **19**(5): 455-477.
- 726 VII. Berlin, K., S. Koren, C. S. Chin, J. P. Drake, J. M. Landolin and A. M. Phillippy (2015). "Assembling  
727 large genomes with single-molecule sequencing and locality-sensitive hashing." Nat Biotechnol **33**(6):  
728 623-630.
- 729 VIII. Briese, T., A. Kapoor, N. Mishra, K. Jain, A. Kumar, O. J. Jabado and W. I. Lipkin (2015). "Virome  
730 Capture Sequencing Enables Sensitive Viral Diagnosis and Comprehensive Virome Analysis." MBio **6**(5):  
731 e01491-01415.
- 732 IX. Brister, J. R., Y. Bao, C. Kuiken, E. J. Lefkowitz, P. Le Mercier, R. Leplae, R. Madupu, R. H.  
733 Scheuermann, S. Schobel, D. Seto, S. Shrivastava, P. Sterk, Q. Zeng, W. Klimke and T. Tatusova (2010).  
734 "Towards Viral Genome Annotation Standards, Report from the 2010 NCBI Annotation Workshop."  
735 Viruses **2**(10): 2258-2268.
- 736 X. Chaisson, M. J. and G. Tesler (2012). "Mapping single molecule sequencing reads using basic  
737 local alignment with successive refinement (BLASR): application and theory." BMC Bioinformatics **13**:  
738 238.
- 739 XI. Chin, C. S., D. H. Alexander, P. Marks, A. A. Klammer, J. Drake, C. Heiner, A. Clum, A. Copeland, J.  
740 Huddleston, E. E. Eichler, S. W. Turner and J. Korlach (2013). "Nonhybrid, finished microbial genome  
741 assemblies from long-read SMRT sequencing data." Nat Methods **10**(6): 563-569.
- 742 XII. Ciuffo, S., S. Kannan, S. Sharma, A. Badretdin, K. Clark, S. Turner, S. Brover, C. L. Schoch, A. Kimchi  
743 and M. DiCuccio (2018). "Using average nucleotide identity to improve taxonomic assignments in  
744 prokaryotic genomes at the NCBI." Int J Syst Evol Microbiol **68**(7): 2386-2392.

- 745XIII. Gargis, A. S., L. Kalman and I. M. Lubin (2016). "Assuring the Quality of Next-Generation  
746 Sequencing in Clinical Microbiology and Public Health Laboratories." *J Clin Microbiol* **54**(12): 2857-2865.
- 747XIV. Goldberg, B., H. Sichtig, C. Geyer, N. Ledebor and G. M. Weinstock (2015). "Making the Leap  
748 from Research Laboratory to Clinic: Challenges and Opportunities for Next-Generation Sequencing in  
749 Infectious Disease Diagnostics." *MBio* **6**(6): e01888-01815.
- 750 XV. Greninger, A. L., K. Messacar, T. Dunnebacke, S. N. Naccache, S. Federman, J. Bouquet, D.  
751 Mirsky, Y. Nomura, S. Yagi, C. Glaser, M. Vollmer, C. A. Press, B. K. Kleinschmidt-DeMasters, S. R.  
752 Dominguez and C. Y. Chiu (2015). "Clinical metagenomic identification of Balamuthia mandrillaris  
753 encephalitis and assembly of the draft genome: the continuing case for reference genome sequencing."  
754 *Genome Med* **7**: 113.
- 755XVI. Hatcher, E. L., S. A. Zhdanov, Y. Bao, O. Blinkova, E. P. Nawrocki, Y. Ostapchuck, A. A. Schaffer  
756 and J. R. Brister (2017). "Virus Variation Resource - improved response to emergent viral outbreaks."  
757 *Nucleic Acids Res* **45**(D1): D482-D490.
- 758XVII. Heger, M. (2018). FDA Mulls Guidelines for NGS-Based Infectious Disease Diagnostics.  
759 GenomeWeb.
- 760XVIII. Kearse, M., R. Moir, A. Wilson, S. Stones-Havas, M. Cheung, S. Sturrock, S. Buxton, A. Cooper, S.  
761 Markowitz, C. Duran, T. Thierer, B. Ashton, P. Meintjes and A. Drummond (2012). "Geneious Basic: an  
762 integrated and extendable desktop software platform for the organization and analysis of sequence  
763 data." *Bioinformatics* **28**(12): 1647-1649.
- 764XIX. Klimke, W., C. O'Donovan, O. White, J. R. Brister, K. Clark, B. Fedorov, I. Mizrachi, K. D. Pruitt and  
765 T. Tatusova (2011). "Solving the Problem: Genome Annotation Standards before the Data Deluge." *Stand*  
766 *Genomic Sci* **5**(1): 168-193.
- 767 XX. Koehler, J. W., A. T. Hall, P. A. Rolfe, A. N. Honko, G. F. Palacios, J. N. Fair, J. J. Muyembe, P.  
768 Mulembekani, R. J. Schoepp, A. Adesokan and T. D. Minogue (2014). "Development and evaluation of a  
769 panel of filovirus sequence capture probes for pathogen detection by next-generation sequencing." *PLoS*  
770 *One* **9**(9): e107007.
- 771XXI. Koren, S., M. C. Schatz, B. P. Walenz, J. Martin, J. T. Howard, G. Ganapathy, Z. Wang, D. A. Rasko,  
772 W. R. McCombie, E. D. Jarvis and M. P. Adam (2012). "Hybrid error correction and de novo assembly of  
773 single-molecule sequencing reads." *Nat Biotechnol* **30**(7): 693-700.
- 774XXII. Koren, S., B. P. Walenz, K. Berlin, J. R. Miller, N. H. Bergman and A. M. Phillippy (2017). "Canu:  
775 scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation." *Genome*  
776 *Res* **27**(5): 722-736.
- 777XIII. Land, M., L. Hauser, S. R. Jun, I. Nookaew, M. R. Leuze, T. H. Ahn, T. Karpinets, O. Lund, G. Kora,  
778 T. Wassenaar, S. Poudel and D. W. Ussery (2015). "Insights from 20 years of bacterial genome  
779 sequencing." *Funct Integr Genomics* **15**(2): 141-161.
- 780XIV. Land, M. L., D. Hyatt, S. R. Jun, G. H. Kora, L. J. Hauser, O. Lukjancenko and D. W. Ussery (2014).  
781 "Quality scores for 32,000 genomes." *Stand Genomic Sci* **9**: 20.
- 782XXV. Langdon, W. B. (2015). "Performance of genetic programming optimised Bowtie2 on genome  
783 comparison and analytic testing (GCAT) benchmarks." *BioData Min* **8**(1): 1.
- 784XXVI. Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin  
785 and S. Genome Project Data Processing (2009). "The Sequence Alignment/Map format and SAMtools."  
786 *Bioinformatics* **25**(16): 2078-2079.
- 787XVII. Morgulis, A., G. Coulouris, Y. Raytselis, T. L. Madden, R. Agarwala and A. A. Schaffer (2008).  
788 "Database indexing for production MegaBLAST searches." *Bioinformatics* **24**(16): 1757-1764.

- ~~789~~VIII. Naccache, S. N., K. S. Peggs, F. M. Mattes, R. Phadke, J. A. Garson, P. Grant, E. Samayoa, S. Federman, S. Miller, M. P. Lunn, V. Gant and C. Y. Chiu (2015). "Diagnosis of neuroinvasive astrovirus infection in an immunocompromised adult with encephalitis by unbiased next-generation sequencing." Clin Infect Dis **60**(6): 919-923.
- ~~793~~XIX. O'Flaherty, B. M., Y. Li, Y. Tao, C. R. Paden, K. Queen, J. Zhang, D. L. Dinwiddie, S. M. Gross, G. P. Schroth and S. Tong (2018). "Comprehensive viral enrichment enables sensitive respiratory virus genomic identification and analysis by next generation sequencing." Genome Res **28**(6): 869-877.
- ~~796~~XX. Roach, D. J., J. N. Burton, C. Lee, B. Stackhouse, S. M. Butler-Wu, B. T. Cookson, J. Shendure and S. J. Salipante (2015). "A Year of Infection in the Intensive Care Unit: Prospective Whole Genome Sequencing of Bacterial Clinical Isolates Reveals Cryptic Transmissions and Novel Microbiota." PLoS Genet **11**(7): e1005413.
- ~~800~~XXI. Schatz, M. C. and B. Langmead (2013). "The DNA Data Deluge: Fast, efficient genome sequencing machines are spewing out more data than geneticists can analyze." IEEE Spectr **50**(7): 26-33.
- ~~802~~XXII. Schlaberg, R., C. Y. Chiu, S. Miller, G. W. Procop, G. Weinstock, C. Professional Practice, M. Committee on Laboratory Practices of the American Society for and P. Microbiology Resource Committee of the College of American (2017). "Validation of Metagenomic Next-Generation Sequencing Tests for Universal Pathogen Detection." Arch Pathol Lab Med.
- ~~806~~XXIII. Simner, P. J., S. Miller and K. C. Carroll (2017). "Understanding the Promises and Hurdles of Metagenomic Next-Generation Sequencing as a Diagnostic Tool for Infectious Diseases." Clin Infect Dis.
- ~~808~~XXIV. Snitkin, E. S., S. Won, A. Pirani, Z. Lapp, R. A. Weinstein, K. Lolans and M. K. Hayden (2017). "Integrated genomic and interfacility patient-transfer data reveal the transmission pathways of multidrug-resistant *Klebsiella pneumoniae* in a regional outbreak." Sci Transl Med **9**(417).
- ~~810~~XXV. Snitkin, E. S., A. M. Zelazny, P. J. Thomas, F. Stock, N. C. S. P. Group, D. K. Henderson, T. N. Palmore and J. A. Segre (2012). "Tracking a hospital outbreak of carbapenem-resistant *Klebsiella pneumoniae* with whole-genome sequencing." Sci Transl Med **4**(148): 148ra116.
- ~~814~~XXVI. Somasekar, S., D. Lee, J. Rule, S. N. Naccache, M. Stone, M. P. Busch, C. Sanders, W. M. Lee and C. Y. Chiu (2017). "Viral Surveillance in Serum Samples From Patients With Acute Liver Failure By Metagenomic Next-Generation Sequencing." Clin Infect Dis.
- ~~817~~XXVII. Stefan, C. P., J. W. Koehler and T. D. Minogue (2016). "Targeted next-generation sequencing for the detection of ciprofloxacin resistance markers using molecular inversion probes." Sci Rep **6**: 25904.
- ~~819~~XXVIII. Tatusova, T., M. DiCuccio, A. Badretdin, V. Chetvernin, E. P. Nawrocki, L. Zaslavsky, A. Lomsadze, K. D. Pruitt, M. Borodovsky and J. Ostell (2016). "NCBI prokaryotic genome annotation pipeline." Nucleic Acids Res **44**(14): 6614-6624.
- ~~822~~XXIX. Tchesnokova, V., M. Billig, S. Chattopadhyay, E. Linardopoulou, P. Aprikian, P. L. Roberts, V. Skrivankova, B. Johnston, A. Gileva, I. Igusheva, A. Toland, K. Riddell, P. Rogers, X. Qin, S. Butler-Wu, B. T. Cookson, F. C. Fang, B. Kahl, L. B. Price, S. J. Weissman, A. Limaye, D. Scholes, J. R. Johnson and E. V. Sokurenko (2013). "Predictive diagnostics for *Escherichia coli* infections based on the clonal association of antimicrobial resistance and clinical outcome." J Clin Microbiol **51**(9): 2991-2999.
- ~~827~~XL. Trombley, A. R., L. Wachter, J. Garrison, V. A. Buckley-Beason, J. Jahrling, L. E. Hensley, R. J. Schoepp, D. A. Norwood, A. Goba, J. N. Fair and D. A. Kulesh (2010). "Comprehensive panel of real-time TaqMan polymerase chain reaction assays for detection and absolute quantification of filoviruses, arenaviruses, and New World hantaviruses." Am J Trop Med Hyg **82**(5): 954-960.

831XLI. Walker, B. J., T. Abeel, T. Shea, M. Priest, A. Abouelliel, S. Sakthikumar, C. A. Cuomo, Q. Zeng, J.  
832 Wortman, S. K. Young and A. M. Earl (2014). "Pilon: an integrated tool for comprehensive microbial  
833 variant detection and genome assembly improvement." PLoS One **9**(11): e112963.  
834XLII. Wilson, M. R., S. N. Naccache, E. Samayoa, M. Biagtan, H. Bashir, G. Yu, S. M. Salamat, S.  
835 Somasekar, S. Federman, S. Miller, R. Sokolic, E. Garabedian, F. Candotti, R. H. Buckley, K. D. Reed, T. L.  
836 Meyer, C. M. Seroogy, R. Galloway, S. L. Henderson, J. E. Gern, J. L. DeRisi and C. Y. Chiu (2014).  
837 "Actionable diagnosis of neuroleptospirosis by next-generation sequencing." N Engl J Med **370**(25):  
838 2408-2417.  
839XLIII. Wilson, M. R., L. L. Zimmermann, E. D. Crawford, H. A. Sample, P. R. Soni, A. N. Baker, L. M. Khan  
840 and J. L. DeRisi (2017). "Acute West Nile Virus Meningoencephalitis Diagnosed Via Metagenomic Deep  
841 Sequencing of Cerebrospinal Fluid in a Renal Transplant Patient." Am J Transplant **17**(3): 803-808.  
842XLIV. Wood, D. E. and S. L. Salzberg (2014). "Kraken: ultrafast metagenomic sequence classification  
843 using exact alignments." Genome Biol **15**(3): R46.  
844

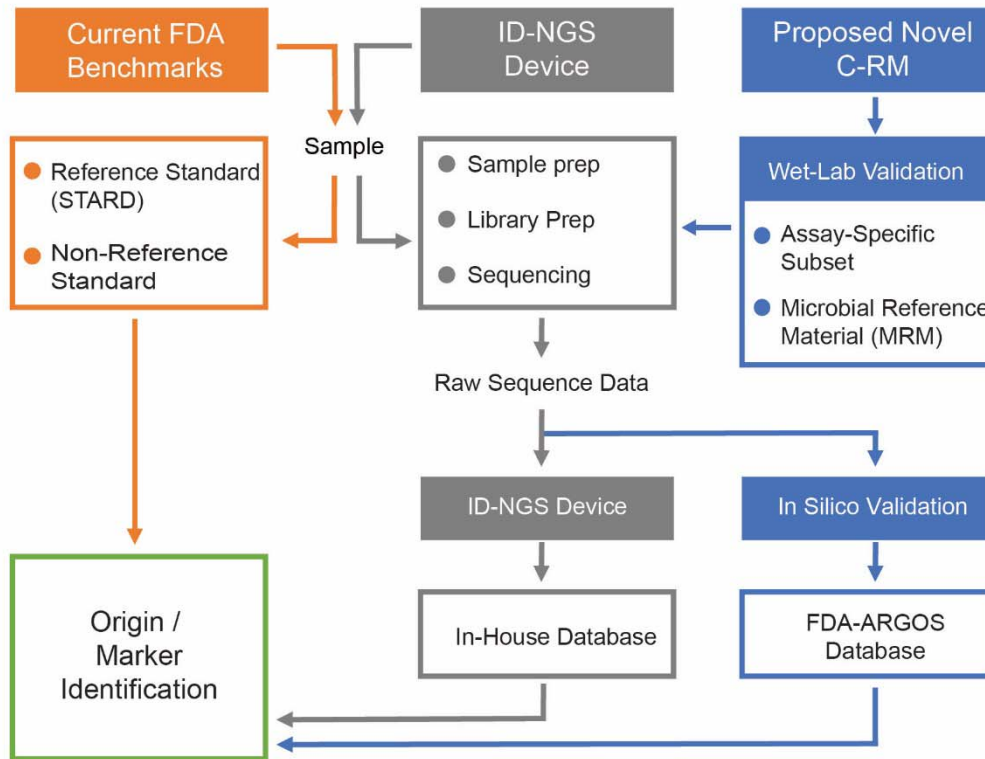
845

846

847 **Figures:**

848 **Figure 1: Proposed Novel Composite Reference Method (C-RM) for ID NGS Diagnostic Assays.**

A.



B.

FDA-ARGOS Quality Metrics	
a.	Organisms identified prior to sequencing by orthogonal reference method
b.	Two sequencing methodologies <ul style="list-style-type: none"> <li>• Bacterial - Long read and short read</li> <li>• Viral - Shotgun, amplicon and RACE</li> </ul>
c.	De-novo assembly
d.	Assembly metrics <ul style="list-style-type: none"> <li>• Coverage : 95% with 20X depth at every position across assembly</li> <li>• N50</li> <li>• L50</li> </ul>
e.	Meets NCBI's taxonomy-specific average nucleotide identity (ANI) thresholds
f.	Meets minimum FDA-ARGOS regulatory grade data requirements

C.

FDA-ARGOS Data Requirements	
a.	Sample Name (Sample ID)
b.	Raw Reads (SRA Accession)
c.	Assemblies (Chromosome, Plasmid, WGS Accession)
d.	Annotations (GenBank Accession)
e.	10-meta data (Biosample Accession)

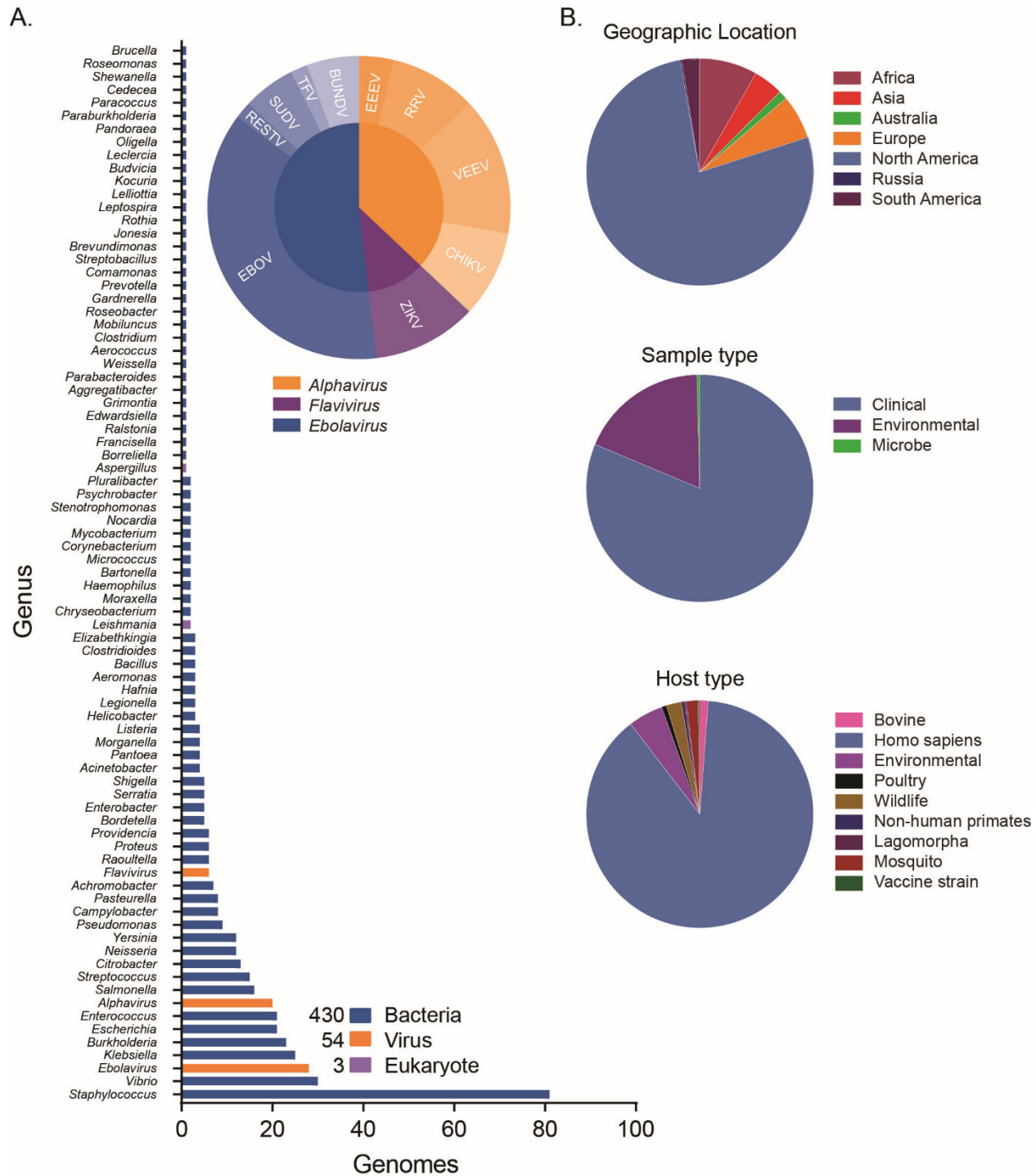
D.

10-Meta Data (Biosample Accession)	
1.	Organism name
2.	Strain name
3.	ID method
4.	Sample type
5.	Host
6.	Isolation provider name
7.	Isolation acquisition ID
8.	NCBI taxonomy ID
9.	Contact name
10.	Clinical or environmental

850 Figure 1A illustrates a walkthrough of the proposed novel composite reference method (C-RM).  
851 Here, we show *in silico* target sequence validation with FDA-ARGOS reference genomes in  
852 combination with a wet lab validation challenge to understand the performance of ID NGS  
853 diagnostic assays. Using raw sequence data from the ID-NGS test device, *in silico* comparison of  
854 results obtained with the assay in-house database to results when using FDA-ARGOS will  
855 evaluate device bioinformatic analysis pipelines and report generation while eliminating the  
856 need for additional sample testing with a gold standard comparator (current FDA benchmarks).  
857 Overall, we anticipate the use of the C-RM method based on assay-specific subsets of clinical  
858 samples and/or microbial reference materials (MRMs) for wet lab validation and FDA-ARGOS  
859 *in silico* target sequence validation to generate scientifically valid evidence for understanding the  
860 performance of ID NGS diagnostic assays. Figure 1B lists the required quality control metrics for  
861 passing the regulatory-grade genome criteria. At a minimum, an FDA-ARGOS regulatory-grade  
862 genome adheres to six metrics (a-f). Specifically, category f details the minimum data  
863 requirements that are further described in Table 1C. In addition, Table 1D lists the 10 critical  
864 meta data that need to be ascribed to a genome to meet the regulatory-grade criteria.  
865

866 **Figure 2: FDA-ARGOS Reference Genome Database.**



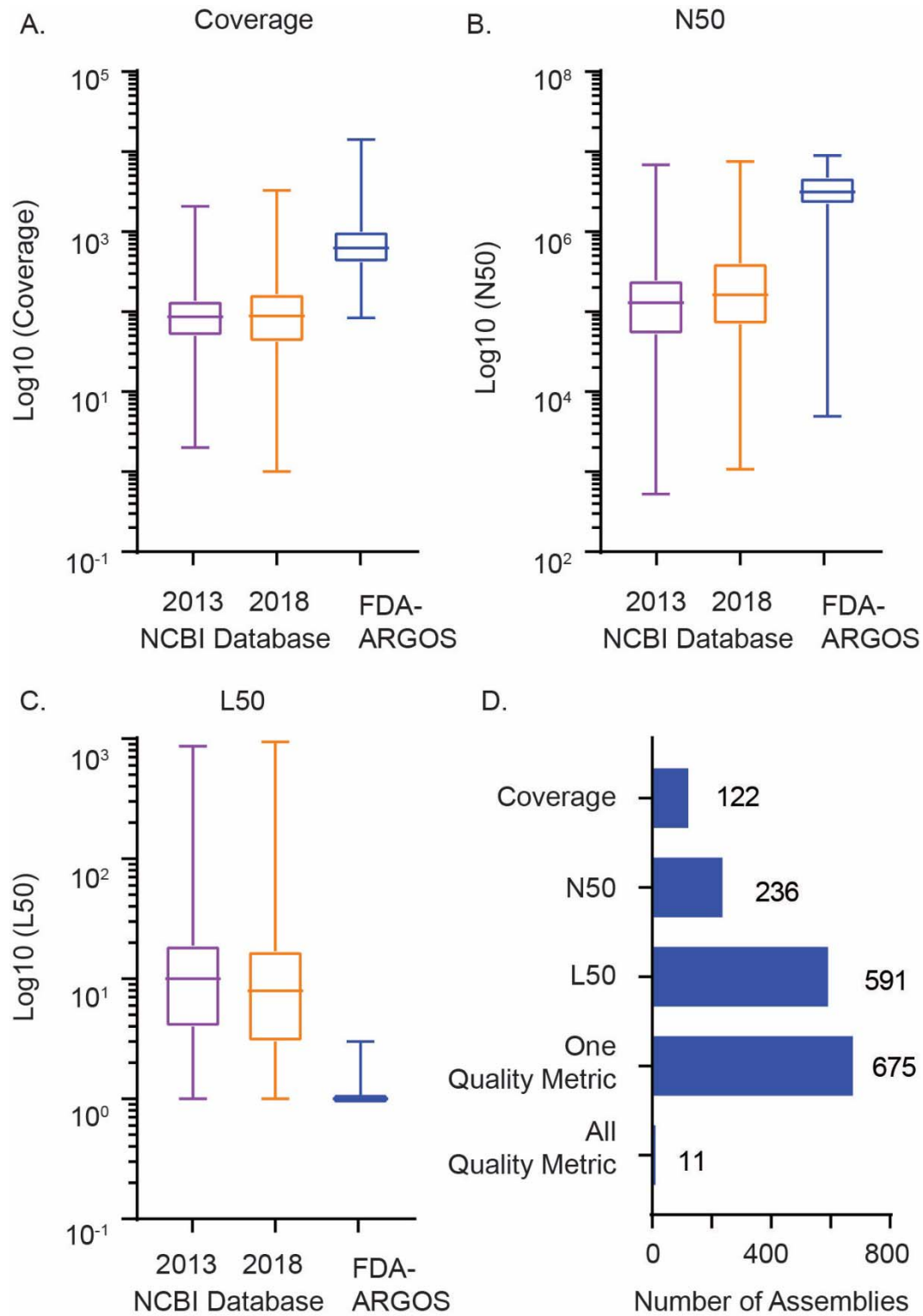


867

868 Summary statistics of the current 487 microbial genomes show primary coverage of FDA-  
 869 ARGOS resides with bacterial isolates, followed by viruses and then eukaryotic parasites (A).  
 870 Supplemental Table 1 provides accessions for all 487 genomes currently available publicly. A

871 majority of FDA-ARGOS constituents (B) originate from North America and are from human  
872 clinical isolation.  
873

874 **Figure 3: FDA-ARGOS Reference Genome Assemblies Quality Metrics.**



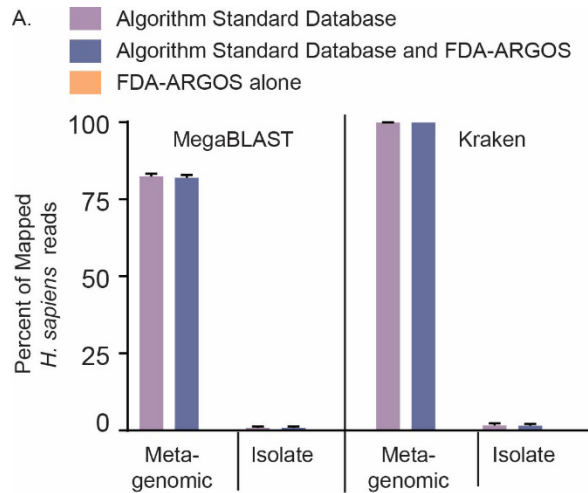
875

876 Comparative microbial genome assembly quality metrics contrasted current FDA-ARGOS  
877 assemblies to 2013 and 2018 NCBI GenBank assemblies submitted for each species captured  
878 within the FDA-ARGOS database. Assembly quality metrics measured included: (A) median  
879 coverage, (B) median N50, (C) median L50 and (D) number of 2018 NCBI genomes that  
880 exhibited all, one or a specific quality control metric used to vet FDA-ARGOS genomes for  
881 inclusion. The NCBI assemblies were downloaded on August 6, 2018.

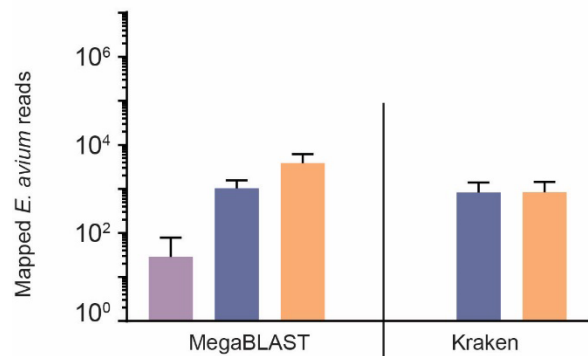
882

883

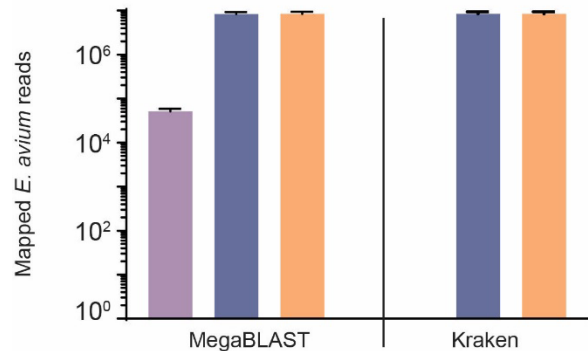
884 **Figure 4 Read Classification Results from Shotgun Sequencing for Identification of**  
885 ***Enterococcus avium*.**



B. Mock clinical *E. avium* reads  
Metagenomics shotgun



C. Clinical *E. avium* reads  
Isolate shotgun



886

887

888

889

890

Visualizing sample analyzed with both MegaBLAST and Kraken percent human reads mapped (A) from metagenomics and isolate sequencing shotgun sequencing showed sequencing clinical matrix resulted in a majority of reads mapping to host background. The number of *E. avium* reads correctly classified from metagenomic samples using three different reference databases

891 (B) varied based on whether MegaBLAST and Kraken standard databases, standard database  
 892 plus FDA-ARGOS, or FDA-ARGOS alone was used for read mapping. Evaluation of *E. avium* reads  
 893 for just the clinical isolate without matrix (C) resulted in a similar relationship of greater number  
 894 of reads mapped when FDA-ARGOS genomes were used in comparison to the algorithm  
 895 standard reference database. All Reference Data Sets and the FDA-ARGOS Database are publicly  
 896 available for data analysis and tool comparison.

897

898 **Tables:**

899

900 **Table 1A: Diagnostic Benchmark and *In Silico* Target Sequence Validation with FDA-ARGOS:**  
 901 **Bundibugyo Performance Summary.**  
 902

Sample	Real-Time PCR (Benchmark)	MIPS (Test Device)	FDA-ARGOS ( <i>In Silico</i> Target Sequence Validation)	
	Quantitation Cycle (Cq) Value	Percent Classified	MegaBLAST Percent Classified	Kraken Percent Classified
2012-1	22.97/22.95	54.95%	59.84%	70.89%
2012-16	ND	0.02%	0.76%	0.02%
2012-91	ND	0.03%	0.65%	0.04%
2012-95	ND	0.02%	0.59%	0.03%
2012-99	ND	0.02%	0.67%	0.06%
2012-120	23.46/23.38	41.87%	45.14%	57.56%
2012-147	25.58/25.52	27.23%	29.78%	47.52%
2012-153	28.14/27.96	38.30%	40.97%	50.70%
2012-176	37.01/36.54	0.01%	0.87%	0.01%
2012-198	ND	0.02%	0.71%	0.03%
NTC	N/A	0.02%	1.59%	1.05%

903 Illustration of an application of the *in silico* target sequence validation method to a targeted ID  
 904 sequencing assay (MIPS). Two bioinformatics tools (MegaBLAST and Kraken) were selected to  
 905 classify reads using default parameters utilizing FDA-ARGOS as a standalone reference  
 906 database. Table 1A showed the traditional benchmark comparison of the MIPS assay to Real-  
 907 Time PCR (RT-PCR) results. Benchmark positive values were only noted for samples that yielded  
 908 duplicative positive results by RT-PCR. Percent reads classified only refer to percentage of reads  
 909 that were assigned to Bundibugyo or Ebola virus, the remaining reads are non-specific or  
 910 "junk".

911

912 **Table 1B: Diagnostic Benchmark and *In Silico* Target Sequence Validation with FDA-ARGOS:**  
 913 **Ebola Makona Performance Summary.**

Sample	Real-Time PCR (Benchmark)	MIPS (Test Device)	FDA-ARGOS ( <i>In Silico</i> Target Sequence Validation)		
	Quantitation Cycle (Cq) Value	Percent Classified	MegaBLAST Percent Classified	Kraken Percent Classified	LMAT Percent Classified
3754-2	35.11/35.72	0.05%	0.60%	0.06%	0.03%
3754-4	33.83/33.36	0.06%	0.68%	0.06%	0.03%
3811-2	36.17/36.14	0.07%	0.56%	0.08%	0.04%
3856-1P	15.95/15.98	76.63%	80.91%	79.50%	42.48%
3913-5	34.00/33.77	0.00%	0.68%	0.00%	0.01%
3958-4	32.77/33.28	0.04%	0.65%	0.05%	0.05%
3991-2	33.92/33.62	0.00%	0.65%	0.00%	0.01%
4007-2	26.30/26.55	21.33%	22.41%	21.81%	12.12%
4015-1	21.66/21.66	74.87%	76.35%	75.82%	39.35%
4033-1	16.59/16.32	76.64%	79.59%	78.97%	41.79%
4268-1P	25.05/25.15	29.71%	30.43%	29.97%	15.87%
4468-3	35.78/35.91	0.04%	0.59%	0.05%	0.03%
4641-3P	31.81/31.82	0.03%	0.59%	0.04%	0.02%
4726-1	21.22/21.21	53.95%	56.44%	54.28%	30.05%
4845-3	35.17/36.71	0.00%	0.66%	0.01%	0.01%
NTC	N/A	0.02%	0.86%	0.04%	0.01%

914 Illustration of an application of the *in silico* target sequence validation method to a targeted ID  
 915 sequencing assay (MIPS). Three bioinformatics tools (MegaBLAST, Kraken and LMAT) were  
 916 selected to classify reads using default parameters utilizing FDA-ARGOS as a standalone  
 917 reference database. Table 1B showed the traditional benchmark comparison of the MIPS assay  
 918 to Real-Time PCR (RT-PCR) results. Benchmark positive values were only noted for samples that  
 919 yielded duplicative positive results by RT-PCR. Percent reads classified only refer to percentage  
 920 of reads that were assigned to Bundibugyo or Ebola virus, the remaining reads are non-specific  
 921 or "junk".

922

923

924 **Table 2: Mock Clinical Evaluation of EBOV NGS Performance.**

925

---

**A. Experimental design and results**

PFU/ml	n	Avg EBOV Reads	Avg %Reads		Positive Samples	Negative Samples
			Mapped	CoV		
1000000 (10X)	16	5442.5	2.66%	136.55%	15	1
500000 (5X)	16	2777.5	2.49%	152.33%	13	3
100000 (1X)	16	351.5	0.58%	247.57%	9	7
NTC	100	4	0.00%	571.69%	1	99

926

### B. Diagnostic performance statistics

N	Positive Predictive Value	Negative Predictive Value	Sensitivity	Specificity	Prevalence
148	97.37% (83.95% to 99.62%)	90.00 % (84.26% to 93.80%)	77.08% (62.69% to 87.97%)	99.00% (94.55% to 99.97%)	32.43% (24.98% to 40.61%)

927

### C. Diagnostic performance statistics for prior probabilities

Prior probability of infection	Positive Predictive Value	Negative Predictive Value
0	0	1
0.01	0.44	1
0.05	0.8	0.99
0.1	0.9	0.97
0.15	0.93	0.96
0.2	0.95	0.95
0.25	0.96	0.93
0.3	0.97	0.91
0.4	0.98	0.87
0.5	0.99	0.81
0.6	0.99	0.74
0.7	0.99	0.65
0.75	1	0.59
0.8	1	0.52
0.85	1	0.43
0.9	1	0.32
0.95	1	0.18
0.99	1	0.04
1	1	0

928 Demonstration (A) of the preliminary diagnostic performance (B) of a targeted ID sequencing  
 929 assay (MIPS) during a mock clinical trial using 48 positive Ebola samples and 100 Ebola negative  
 930 samples. Numbers in parentheses represent the 95% Confidence Interval. Positive and negative  
 931 predictive values are shown for prior probabilities (C) of infection ranging from 0-1.