1    **Rapid metagenomic identification of viral pathogens in clinical samples by**

2    **real-time nanopore sequencing analysis**

3

4    Alexander L. Greninger[1,2], Samia N. Naccache[1,2,#], Scot Federman[1,2,#], Guixia

5    Yu[1,2], Placide Mbala[3,6], Vanessa Bres[4], Doug Stryke[1,2], Jerome Bouquet[1,2],

6    Sneha Somasekar[1,2], Jeffrey M. Linnen[4], Roger Dodd[5], Prime Mulembakani[6],

7    Bradley S. Schneider[6], Jean-Jacques Muyembe[3], Susan L. Stramer[5], Charles Y.

8    Chiu[1,2]

9

10   Affiliations:

11   [1]Department of Laboratory Medicine, University of California, San Francisco, CA,

12   USA 94107

13   [2]UCSF-Abbott Viral Diagnostics and Discovery Center, San Francisco, CA, USA,

14   91407

15   [4]Institut National de Recherche Biomédicale, Democratic Republic of the Congo,

16   Kinshasa, Africa

17   [5]Hologic, Inc., Bedford, MA 01730

18   [5]American Red Cross, Gaithersburg, MD  2087

19   [6]Metabiota, Inc., San Francisco, CA  94104

20   [7]Department of Medicine, Division of Infectious Diseases, University of California,

21   San Francisco

22

23   [#]These authors contributed equally to this work.

24

25

**ABSTRACT**

We report unbiased metagenomic detection of chikungunya virus (CHIKV), Ebola

virus (EBOV), and hepatitis C virus (HCV) from four human blood samples by

MinION nanopore sequencing coupled to a newly developed, web-based pipeline

for real-time bioinformatics analysis on a computational server or

laptop(MetaPORE). At titers ranging from $10^7$-$10^8$ copies per milliliter, reads to

EBOV from two patients with acute hemorrhagic fever and CHIKV from an

asymptomatic blood donor were detected within 4 to 10 minutes of data

acquisition, while lower titer HCV virus ($1 \times 10^5$ copies per milliliter) was detected

within 40 minutes. Analysis of mapped nanopore reads alone, despite an

average individual error rate of 24% [range 8-49%], permitted identification of the

correct viral strain in all 4 isolates, and 90% of the genome of CHIKV was

recovered with >98% accuracy. Using nanopore sequencing, metagenomic

detection of viral pathogens directly from clinical samples was performed within

an unprecedented <6 hours sample-to-answer turnaround time and in a

timeframe amenable for actionable clinical and public health diagnostics.

**Background**

Acute febrile illness has a broad differential diagnosis and can be caused by a variety of pathogens. Metagenomic next-generation sequencing (NGS) is particularly attractive for diagnosis and public health surveillance of febrile illness because the approach can broadly detect viruses, bacteria, and parasites in clinical samples by uniquely identifying sequence data[1, 2]. Although currently limited by a sample-to-answer turnaround time of >20 hr (Fig. 1), we and others have reported that unbiased pathogen detection using metagenomic NGS can generate actionable results in timeframes relevant to clinical diagnostics [3-6] and public health [7, 8]. As sequence reads are generated in parallel and not in series, real-time analysis by second-generation platforms such as Illumina and Ion Torrent has been hampered by the need to wait until a sufficient read length has been achieved for diagnostic pathogen identification.

Nanopore sequencing is a third-generation sequencing technology that has two key advantages over second-generation technologies – longer reads and the ability to perform real-time sequence analysis. To date, the longer nanopore reads have enabled scaffolding of prokaryotic and eukaryotic genomes and sequencing of bacterial and viral cultured isolates[9-12], but the platform's capacity for real-time metagenomic analysis of clinical samples has not yet been leveraged. As of mid-2015, the MinION nanopore sequencer is capable of producing at least 100,000 sequences with an average read length of 5 kB, in total producing up to 1 Gb of sequence in 24 hours on one flow cell [13]. Here we present using nanopore sequencing for metagenomic detection of viral

65  pathogens from clinical samples with sample-to-answer turnaround times of

66  under 6 hours.  We also present MetaPORE, a real-time, web-based sequence

67  analysis and visualization tool for pathogen identification from nanopore data.

68

69  **Materials and Methods**

70

71  **MAP program**

72      Since July 2014, our lab has participated in the MinION Access Program

73  (MAP), an early access program for beta users of the Oxford Nanopore MinION.

74  Program participants receive free flow cells and library preparation kits for testing

75  and validation of new protocols and applications on the MinION platform.  During

76  our time in the MAP program, we have seen significant progress in quality control

77  and sequencing yield of flow cells (Table 1).

78

79  **Nucleic acid extraction**

80      Total nucleic acid was extracted from 400 μL of a chikungunya virus

81  (CHIKV)-positive serum sample previously collected from an asymptomatic blood

82  donor during the 2014 CHIKV outbreak in Puerto Rico (Chik1) [14].  The serum

83  sample was inactivated in a 1:3 ratio of TRIzol LS (Life Technologies, Carlsbad,

84  CA, USA) at the American Red Cross prior to shipping to University of California,

85  San Francisco (UCSF).  Direct-zol RNA MiniPrep (Zymo Research, Irvine, CA,

86  USA) was used for nucleic acid extraction, including on-column treatment with

87    Turbo DNAse (Life Technologies) for 30 min at 37°C to deplete human host

88    genomic DNA.

89        For the Ebola virus (EBOV) specimens, total nucleic acid was extracted

90    using a QIAamp Viral RNA kit (Qiagen, Valencia, CA, USA) from 140 $\mu$L of whole

91    blood from two patients with suspected Ebola hemorrhagic fever from a 2014

92    outbreak in the Democratic Republic of the Congo (DRC) (Ebola1 and Ebola2).

93    RNA was extracted at Institut National de Recherche Biomédicale in Kinshasa,

94    DRC, preserved using RNAstable (Biomatrica, San Diego, CA, USA), and

95    shipped at room temperature to UCSF.  Upon receipt, the extracted RNA sample

96    was treated with 1 $\mu$L Turbo DNase (Life Technologies), followed by clean-up

97    using the Direct-zol RNA MiniPrep Kit (Zymo Research).

98        For the HCV sample, an HCV-positive serum specimen at a titer of

99    $1.6 \times 10^7$ cp/mL (HepC1) was diluted to $1 \times 10^5$ cp/mL using pooled negative serum.

100    Total nucleic acid was then extracted from 400uL of serum using the EZ1 Viral

101    RNA kit followed by treatment with Turbo DNase at 30 min at 37°C and clean-up

102    using the RNA Clean and Concentrator Kit (Zymo Research).

103

104    **Molecular confirmation of viral infection**

105        A previously reported TaqMan quantitative reverse-transcription PCR

106    (qRT-PCR) assay targeting the EBOV NP gene was used for detection of EBOV

107    and determination of viral load.  The assay was run on a Stratagene MX300P

108    real-time PCR instrument and performed using the TaqMan Fast Virus 1-Step

109    Master Mix (Life Technologies) in 20 µL total reaction volume (5 µL 4x Taqman

110    mix, 1ul sample extract), with 0.75uM of each primer (F565 5'-

111    TCTGACATGGATTACCACAAGATC-3', R640 5'-

112    GGATGACTCTTTGCCGAACAATC-3') and 0.6uM of the probe (p597S 6FAM-

113    AGGTCTGTCCGTTCAA-MGBNFQ). Conditions for the qRT-PCR were as

114    follows: 50°C, 10minutes / 95°C, 20 s followed by 45 cycles of 95°C 3s / 60°C, 30

115    s. Viral copy number was calculated by standard curve analysis with a plasmid

116    vector containing the EBOV amplicon. The first EBOV sample analyzed by

117    nanopore sequencing (Ebola1) corresponded to the Ebola virus/*H.sapiens*-

118    wt/COD/2014/Lomela-Lokolia16 strain, while the second Ebola sample (Ebola2)

119    corresponded to the Ebola virus/*H.sapiens*-wt/COD/2014/Lomela-LokoliaB11

120    strain. The CHIKV-positive sample was identified and quantified using a

121    transcription-mediated amplification (TMA) assay (Hologic, Bedford, MA, USA)

122    as previously described [14]. HCV was quantified using the FDA-approved

123    Abbott RealTi*me* RT-PCR assay as performed in the UCSF Clinical Microbiology

124    Laboratory on the Abbott Molecular m2000 system.

125

126    **Construction of metagenomic amplified cDNA libraries**

127        To obtain $\geq$1 $\mu$g of metagenomic cDNA library required for the nanopore

128    sequencing protocol, randomly amplified cDNA was generated using a primer-

129    extension pre-amplification method (Round A/B) as described previously[15-17].

130    Of note, this protocol has been extensively tested on clinical samples for

131    metagenomic pan-pathogen detection of DNA and RNA viruses, bacteria, fungi,

132    and parasites [4, 6, 15, 17, 18]. Briefly, in Round A, RNA was reverse-

133  transcribed with SuperScript III Reverse Transcriptase (Life Technologies,) using

134  Sol-PrimerA (5'-GTTTCCCACTGGAGGATA-$N_9$-3'), followed by second-strand

135  DNA synthesis with Sequenase DNA polymerase (Affymetrix, Santa Clara, CA,

136  USA).  Reaction conditions for Round A were as follows: 1 µL of Sol-PrimerA (40

137  pmol/µl) was added to 4 µl of sample RNA, heated at 65°C for 5 minutes, then

138  cooled at room temperature for 5 minutes.  Five µL of SuperScript Master Mix (2

139  µl 5X First-Strand Buffer, 1 µL water, 1 µL 12.5 mM dNTP mix, 0.5 µL 0.1M DTT,

140  0.5 µL SS III RT) was then added and incubated at 42°C for 60 minutes.  For

141  second strand synthesis, 5 µL of Sequenase Mix #1 (1 µL 5X Sequenase Buffer,

142  3.85 µL ddH$_2$O, 0.15  µL Sequenase enzyme) was added to the reaction mix and

143  incubated at 37°C x 8 min, followed by addition of Sequenase Mix #2 (0.45 µl

144  Sequenase Dilution Buffer, 0.15 µl Sequenase Enzyme) and a second incubation

145  at 37°C x 8 min.  Round B reaction conditions were as follows: 5 µL of Round A-

146  labeled cDNA was added to 45 µL of KlenTaq master mix per sample (5 µL 10X

147  KlenTaq PCR buffer, 1 µL 12.5 mM dNTP, 1 µL 100 pmol/µl Sol-PrimerB (5'-

148  GTTTCCCACTGGAGGATA-3'), 1 µL KlenTaq LA (Sigma-Aldrich, St Louis, MO),

149  37 µL ddH$_2$O).  Reaction conditions for the PCR were as follows: 94°C for 2

150  minutes; 25 cycles of 94°C for 30 sec, 50°C for 45 sec, 72°C for 60 sec; 72°C for

151  5 minutes.

152

153  **Preparation of nanopore sequencing libraries**

154      Amplified cDNA from Round B was purified using AMPure XP beads

155  (Beckman Coulter, Brea, CA), and 1 µg DNA was used as input into Oxford

156 Nanopore Genomic DNA MAP-003 Kits (Chik1, Ebola1) and MAP-004 Kits

157 (HepC1, Ebola2) for generation of MinION Oxford Nanopore-compatible

158 libraries[9, 11]. Briefly, the steps include (1) addition of control lambda phage

159 DNA, (2) end-repair with the NEBNext End Repair Module, (3) 1X AMPure

160 purification, (4) dA-tailing with the NEBNext dA-tailing Module, (5) ligation to

161 protein-linked adapters HP/AMP (Oxford Nanopore Technologies, Oxford, UK)

162 using the NEBNext QuickLigation Module x 10 min at room temperature, (6)

163 purification of ligated libraries using magnetic His-Tag Dynabeads (Life

164 Technologies), and (7) elution in 25 µL buffer (Oxford Nanopore Technologies).

165 Lambda phage DNA was not added during preparation of the Ebola2 sample

166 library.

167

**Nanopore sequencing**

169 Nanopore libraries were run on an Oxford Nanopore MinION flow cell

170 after loading 150 µL sequencing mix (6 µL library, 3 µL fuel mix, 141 µL BP mix)

171 per the manufacturer's instructions. The Chik1 and Ebola1 samples were run

172 consecutively on the same flow cell, with an interim wash performed using Wash-

173 Kit-001 (Oxford Nanopore).

174

**Illumina sequencing**

176 For the Chik1 and Ebola1 specimens, amplified round B cDNA was

177 purified using AMPure XP beads (Beckman Coulter) and 2 ng used as input into

178 the Nextera XT Kit (Illumina). After 13 cycles of amplification, Illumina library

179 concentration and average fragment size were determined using the Agilent

180 Bioanalyzer. Sequencing was performed on an Illumina MiSeq using 150

181 nucleotide (nt) single-end runs and analyzed using the SURPI computational

182 pipeline (UCSF, CA, USA) [15].

183

184 **MetaPORE bioinformatics pipeline**

185 We developed a custom bioinformatics pipeline for real-time pathogen

186 identification and visualization from nanopore sequencing data (MetaPORE),

187 available at https://github.com/chiulab/MetaPORE. The MetaPORE pipeline

188 consists of a set of Linux shell scripts, Python programs, and JavaScript / HTML

189 code, and was tested and run on an Ubuntu 14.10 computational server with 64

190 cores and 512 GB memory. Alternatively, MetaPORE was tested and run on a

191 laptop (4 cores, 32GB RAM, Ubuntu 14.10) by restricting the identification

192 database to viral sequences, rather than all of NCBI nt.

193 Raw FAST5/HDF files from the MinION instrument are base-called using

194 the Metrichor 2D Basecalling v1.14 pipeline (Metrichor). The MetaPORE pipeline

195 continually scans the Metrichor download directory for batch analysis of

196 downloaded sequence reads. For each batch of files (collected every time 200

197 reads are downloaded in the download directory or ≥2 min of elapsed time), the

198 2D read or either the template or complement read, depending on which is of

199 higher quality, is converted into a FASTQ file using HDF5 Tools

200 (https://www.hdfgroup.org/HDF5/doc/RM/Tools.html). The *cutadapt* program is

201 then used to trim Sol-PrimerB adapter sequences from the ends of the reads[19].

202   Next, the nucleotide BLAST (BLASTn) aligner is used to computationally subtract

203   host reads [15, 20] aligning to the human fraction of the National Center for

204   Biotechnology Information (NCBI) nucleotide collection database (NT database,

205   downloaded March 2015) at word size 11 and e-value cutoff of $10^{-5}$.  On our 64

206   core machine, the remaining, non-human reads are then aligned by BLASTn to

207   the entire NCBI NT database using the same parameters.  On a laptop, the non-

208   human reads are aligned to the viral fraction of the NCBI NT database. Reads

209   that hit this viral database are then aligned by BLASTn to NCBI nt. For each

210   read, the single best hit by e-value is retained.  The NCBI GenBank gene

211   identifier assigned to the best hit is then annotated by taxonomic lookup of the

212   corresponding lineage, family, genus, and species [15].

213        For real-time visualization of results, a graphical user interface was

214   developed for the MetaPORE pipeline.  A "live" taxonomic count table is

215   displayed as a donut chart using the CanvasJS (http://canvasjs.com) graphics

216   suite, with the chart refreshing every 30 sec (Supplementary Data, Movie 1).  For

217   each viral family, genus, and species, the "top hit" is chosen to be the reference

218   sequence with the greatest number of aligned reads, with priority given to

219   reference sequences in the following order: (1) complete genomes; (2) complete

220   sequence, or (3) partial sequences / individual genes.  Coverage maps are

221   generated in MetaPORE by mapping all aligned reads at a given taxonomic level

222   (species, genus, or family) to the 'top hit" using LASTZ v1.02 [21], with interactive

223   visualization provided using a custom web program that accesses the

224   HighCharts JavaScript library (http://www.highcharts.com).  A corresponding

225    interactive pairwise identity plot is generated by using SAMtools [22] to calculate

226    the consensus FASTA sequence from the coverage map, followed by pairwise

227    100 bp sliding-window comparisons of the consensus to the reference sequence

228    using the BioPython implementation of the Needleman-Wunsch algorithm [23,

229    24]. For purposes of comparison, the MetaPORE pipeline was also run on a

230    subset of 100,000 reads from parallel Illumina MiSeq data corresponding to the

231    Chik1, Ebola1, and Ebola2 samples.

232

233    **Phylogenetic Analysis**

234      The overall CHIKV phylogeny (Fig. 2E, inset) consisted of all 188 near-

235    complete or complete genome CHIKV sequences available in the NCBI NT

236    database as of March 2015. A subphylogeny including the MiSeq- and

237    nanopore-sequenced Puerto Rico strain PR-S6 presented here and

238    previously[14], as well as additional Caribbean CHIKV strains and other

239    representative members of the Asian-Pacific clade, was also analyzed (Fig. 2E).

240    The EBOV phylogeny (Fig. 3E) consisted of the newly MiSeq- and nanopore-

241    sequenced Ebola strain Lomela-LokoliaB11 from the 2014 DRC outbreak [25]

242    and other representative EBOV strains, including strains from the 2014-2015

243    West African outbreak[8, 26]. Sequences were aligned using the MAFFT

244    algorithm[27], and phylogenetic trees were constructed using the MrBayes

245    algorithm [28] in the Geneious software package [29].

246

247

## RESULTS

### Example 1: High-titer CHIKV (Flowcell #1)

To test the ability of nanopore sequencing to identify metagenomic reads from a clinical sample, we first analyzed a plasma sample harboring high-titer CHIKV and previously sequenced on an Illumina MiSeq platform (Fig. 2B) [14].  The plasma sample corresponded to an asymptomatic blood donor who had screened positive for CHIKV infection during the 2014 outbreak in Puerto Rico (strain PR-S6), with a calculated viral titer of $9.1x10^7$ copies/mL.

A read aligning to CHIKV, the 96th read, was sequenced within 6 min (Fig. 2B) and detected by BLASTN alignment to the NCBI NT database within 8 min of data acquisition, demonstrating an overall sample-to-detection turnaround time of <6 hr (Fig. 1).  After early termination of the sequencing run at the 2 hr 15 min time point, 556 of 19,452 total reads (2.85%) were found to align to CHIKV (Fig. 2B and C).  The individual CHIKV nanopore reads had an average length of 455 nucleotides (nt) [range 126-1477] and average percent identity of 79% [range 51-92%] to the most closely matched reference strain, a CHIKV strain from the neighboring British Virgin Islands (KJ451624) (Table 1).  When only high-quality 2D pass reads were included, 346 of 5,139 (6.7%) reads aligned to CHIKV, comparable to the proportion of CHIKV reads (248,677 of 3,235,096, 7.7%) identified by corresponding metagenomic sequencing on the Illumina MiSeq (Fig. S1).

Mapping of the 556 nanopore reads aligning to CHIKV to the assigned reference genome (KJ451624) showed recovery of 90% of the genome at 3X

271  coverage and 98% at 1X coverage.  Notably, despite high individual mapped

272  error rates in the nanopore reads, ranging from 8 to 49%, 97-99% identity to the

273  reference genome (KJ451624) was achieved across contiguous regions with at

274  least 3X coverage (Fig. 2D).  Furthermore, phylogenetic analysis revealed co-

275  clustering of the CHIKV genomes independently assembled from MinION

276  nanopore and Illumina MiSeq reads on the same branch within the Caribbean

277  subclade (Fig. 2E).  Overall, a large proportion of reads (55%) in the error-prone

278  nanopore data remained unidentifiable, while other aligning reads aside from

279  CHIKV corresponded to human, lambda phage control spike-in, uncultured

280  bacterial, or other eukaryotic sequences (Fig. 2C).

281

282  **Example 2: High-titer Ebola virus (Flowcell #1)**

283          We next attempted to replicate our metagenomic detection result on the

284  nanopore sequencer with a different virus by testing a whole-blood sample from

285  a patient with Ebola hemorrhagic fever during the August 2014 outbreak in the

286  DRC (Ebola1, strain Lomela-Lokolia16) [25].  To conserve flowcells, the same

287  nanopore flowcell used to run the Chik1 sample was washed and stored

288  overnight at 4°C, followed by nanopore sequencing of the Ebola1 sample (viral

289  titer of $1.0 \times 10^7$ copies/mL by real-time qRT-PCR).   Only 41 of 13,090 nanopore

290  reads (0.31%) aligned to EBOV, as compared to 20,820 of 2,743,589 reads

291  (0.76%) at 117X coverage for the Illumina MiSeq (Fig. 3B and D; Fig S1).  The

292  decrease in relative number (41 versus 556) and percentage (0.31% versus

293  2.85%) of target viral nanopore reads in the Ebola1 relative to Chik1 sample is

294   consistent with the lower levels of viremia ($1.0 \times 10^7$ versus $9.1 \times 10^7$ copies/mL)

295   and high host background (whole blood versus plasma) (Fig. 2C).  Nonetheless,

296   the first read aligning to EBOV was detected in a similar timeframe as in the

297   Chik1 sample, sequenced within 8 min and detected within 10 min of data

298   acquisition.  EBOV nanopore reads were 359 nt in length on average [range 220-

299   672 nt], with a mean individual pairwise identity of 78% [range 56-88%].

300   Nevertheless, the majority of Ebola sequences (31 of 41, 76%) were found to

301   align to the most closely matched strain of EBOV in the NT reference database,

302   strain Lomela-Lokolia16 (Fig. 2D).

303       Despite washing the flowcell between the two successive runs, 7 CHIKV

304   reads were recovered during the Ebola1 library sequencing, suggesting the

305   potential for carryover contamination.  CHIKV reads were not present in the

306   corresponding Illumina MiSeq run (Fig. S1-B), confirming that the source of the

307   contamination derived from the Chik1 library that was run on the same flow cell

308   as the Ebola1 library.

309

310

311   **Example 3: Moderate-titer HCV (Flowcell #2)**

312       Our previous experiments revealed both the total number of

313   metagenomic reads and proportion of target viral reads at a given titer that could

314   be obtained from a single MinION flowcell, and showed that the proportion of

315   viral reads obtained by metagenomic nanopore and MiSeq sequencing was

316   comparable. Thus, we projected that the minimum concentration of virus that

317   could be reproducibly detected using our current metagenomic protocol would be

318   $1x10^5$ copies/mL.  An HCV-positive clinical sample (HepC1) was diluted in

319   negative control serum matrix to a titer of $1x10^5$ copies/mL and processed for

320   nanopore sequencing using an upgraded library preparation kit (MAP-004).  After

321   four consecutive runs on the same flowcell with repeat loading of the same

322   metagenomic HepC1 library (Fig. 3A), a total of 85,647 reads were generated, of

323   which only 6 (0.0070%) aligned to HCV (Fig. 3B).  Although the entire series of

324   flowcell runs lasted for >12 hr, the first HCV read was sequenced within 34 min

325   (Fig. 3B), enabling BLASTn detection within 36 min of data acquisition. Given the

326   low titer of HCV in the HepC1 sample and hence low corresponding fraction of

327   HCV reads in the nanopore data, the vast majority (96%) of viral sequences

328   identified corresponded to the background lambda phage spike-in.  Importantly,

329   although only 6 HCV reads were identified by nanopore sequencing, all 6 reads

330   aligned to the correct genotype, genotype 1b (Fig. 3D).

331

332   **Example 4: High-titer Ebola virus with real-time MetaPORE analysis**

333   **(Flowcell #3)**

334        To enable real-time analysis of nanopore sequencing data, we

335   combined our BLASTn analysis with monitoring and user-friendly web

336   visualization into a real-time bioinformatics pipeline for pathogen detection

337   named MetaPORE.  We tested MetaPORE by sequencing a nanopore library

338   (Ebola2) constructed using the upgraded MAP-004 kit and corresponding to a

339   whole blood sample from a patient with suspected Ebola hemorrhagic fever

340 during the 2014 DRC outbreak. Four consecutive runs of the Ebola2 library on

341 the same flowcell over 34 hr, yielded a total of 335,044 reads, of which 593

342 (0.18%) aligned to EBOV (141 of 6009, or 2.3%, of 2D pass reads). Notably, the

343 first EBOV read was sequenced 44 sec after data acquisition and correctly

344 detected in ~3 min by MetaPORE (Fig. 4B; Supplementary Information: Movie).

345 A total of 3 EBOV reads, mapping across the genome and confirming

346 unambiguous detection of the virus, were identified using MetaPORE within 9

347 min of data acquisition (Supplementary Information: Movie). In the corresponding

348 Illumina MiSeq data, 1,778 reads (0.93%) out of a total of 192,158 reads aligned

349 to EBOV by BLASTn, comparable to the proportion of EBOV reads in nanopore

350 2D pass data. Nanopore read coverage across the EBOV genome was

351 relatively uniform with at least 1 read mapping to >88% of the genome and areas

352 of zero coverage also seen with much higher-coverage Illumina MiSeq data (Fig.

353 4D). The detection of EBOV by real-time metagenomic nanopore sequencing

354 was confirmed by qRT-PCR testing of the clinical blood sample, which was

355 positive for EBOV at an estimated titer of 7.64 x$10^7$ cp/ml. Phylogenetic analysis

356 of the Ebola2 genome independently recovered by MinION nanopore and

357 Illumina MiSeq sequencing revealed that nanopore sequencing alone was

358 capable of pinpointing the correct EBOV outbreak strain and country of origin

359 (Fig. 4E).

360

361

362

**DISCUSSION**

363

364    Unbiased point-of-care testing for pathogens by rapid metagenomic

365    sequencing has the potential to radically transform infectious disease diagnosis

366    in clinical and public health settings.  In this study, we sought to demonstrate the

367    potential of the nanopore instrument for metagenomic pathogen identification in

368    clinical samples by coupling an established assay protocol with a new, real-time

369    sequence analysis pipeline. To date, high reported error rates (10-30%) and

370    relatively low-throughput (<100,000 reads per flow cell) have hindered the utility

371    of nanopore sequencing for analysis of metagenomic clinical samples [9, 11].

372    Prior work on infectious disease diagnostics using nanopore has focused on

373    rapid PCR amplicon sequencing of viruses and bacteria[11], or real-time

374    sequencing of pure bacterial isolates in culture, such *Salmonella* in a hospital

375    outbreak [12].  To our knowledge, this is the first time that nanopore sequencing

376    has been used for real-time metagenomic detection of pathogens in complex,

377    high-background clinical samples in the setting of human infections.  Here we

378    also sequenced a microbial genome to high accuracy (>98% identity) directly

379    from a metagenomic clinical sample and not from culture, using only multiple

380    overlapping, albeit error-prone, nanopore reads and without resorting to the use

381    of a secondary platform such as an Illumina MiSeq for sequence correction (Fig.

382    2D).

383    Real-time analysis is critical for time-critical sequencing applications such

384    as outbreak investigation[7] and metagenomic diagnosis of life-threatening

385    infections in hospitalized patients[3, 4, 6].  NGS analysis for clinical diagnostics is

386 currently performed after sequencing is completed, analogous to how PCR

387 products were analyzed by agarose gel electrophoresis in the 1990s. To date,

388 most clinical PCR assays have been converted to a real-time format that reduces

389 "hands-on" laboratory technician time and effort, decreases overall sample-to-

390 answer turnaround times, and potentially yields quantitative information. Notably,

391 our nanopore data suggest that very few reads are needed to provide an

392 unambiguous diagnostic identification, despite high individual per read error rates

393 of 10-30%. The ability of nanopore sequence analysis to accurate identify

394 viruses to the species and even strain / genotype level is facilitated by the high

395 specificity of sequence data, especially with the longer target viral reads

396 achievable by nanopore (Table 1, 391 bp average length) versus second-

397 generation sequencing.

398 Although the overall turnaround time from metagenomic sample-to-

399 detection has been reduced to <6 hr, many challenges remain for routine

400 implementation of this technology in clinical and public health settings.

401 Improvements to make library preparation faster and more robust are critical,

402 including automation and optimization of each step in the protocol. We also

403 looked only at clinical samples at moderate to high titers of $10^5$-$10^8$ copies / mL,

404 and the sensitivity of metagenomic nanopore sequencing at lower titers remains

405 unclear at current achievable sequencing depths. Standard wash protocols

406 appear inadequate to prevent cross-contamination when reusing the same flow

407 cell, as CHIKV reads were identified in the downstream Ebola1 sample sequence

408 run. One solution may be to perform only one nanopore sequencing run per flow

409 cell for clinical diagnostic purposes, akin to how disposable cartridges are used

410 for clinical quantitative PCR testing on a Cepheid GenXpert instrument to prevent

411 cross-contamination[30]. Another is to uniquely barcode individual samples at

412 the cost of added time and effort. Finally, the current accuracy of a single

413 nanopore read will most likely be insufficient to allow confident species

414 identification of bacteria, fungi, or parasites, which have much larger genomes

415 and shared conserved genomic regions than viruses. Single-nucleotide

416 resolution will also be required for detection of antimicrobial resistance markers

417 [31], which is difficult to achieve from relatively low-coverage metagenomic data

418 [32]. These limitations can potentially be overcome by target enrichment

419 methods such as capture probes to increase coverage, improvements in

420 nanopore sequencing technology, or more accurate base-calling and alignment

421 algorithms for nanopore data [33, 34].

422

423

424


425 **CONCLUSION**

426 Our results indicate that unbiased metagenomic detection of viral

427 pathogens from clinical samples with a sample-to-answer turnaround time of <6

428 hours and real-time bioinformatics analysis is feasible with nanopore sequencing.

429 We demonstrate unbiased detection of diagnostic EBOV sequence in under four

430 minutes after the Oxford Minion nanopore initiated sequencing. This technology

431 will be particularly desirable for enabling point-of-care genomic analyses in the

432    developing world, where critical resources, including reliable electric power,

433    laboratory space, and computational server capacity, are often severely limited.

434    MetaPORE, the real-time sequencing analysis platform developed here, is web-

435    based and able to be run on a laptop.  As sequencing yield, quality, and

436    turnaround times continue to improve, we anticipate that third-generation

437    technologies such as nanopore sequencing will challenge clinical diagnostic

438    mainstays such as PCR and TMA testing, fulfilling the dream of an unbiased,

439    point-of-care test for infectious diseases.

440

441    **FIGURE LEGENDS**

442

443    **Figure 1.  Metagenomic sequencing workflow for MinION nanopore**

444    **sequencing compared to Illumina MiSeq sequencing.**  The turnaround time

445    for sample-to-detection nanopore sequencing, defined here as the cumulative

446    time taken for nucleic acid extraction, reverse transcription, library preparation,

447    sequencing, MetaPORE bioinformatics analysis, and detection, was under 6 hr,

448    while Illumina sequencing took over 20 hr.  The time differential is accounted for

449    by increased times for library quantitation, sequencing, and bioinformatics

450    analysis with the Illumina protocol.  *assumes a 12-hr 50 nucleotide (nt) single-

451    end MiSeq run of ~12-15 million reads, with 50 nt the minimum read length

452    needed for accurate pathogen identification. **denotes estimated average SURPI

453    bioinformatics analysis run length for MiSeq data[15].  The stopwatch is depicted

454    as a 12-hr clock.

455

**Figure 2. Metagenomic identification of CHIKV and EBOV from clinical blood samples by nanopore sequencing. (A)** Timeline of sequencing runs on flowcell #1 with sample reloading, plotted as a function of elapsed time in hr since the start of flowcell sequencing. **(B)** Cumulative numbers of all sequenced reads (black line) and target viral reads (red line) from the Chik1 run (left) or Ebola1 run (right ), plotted as a function of individual sequencing run time in min. **(C)** Taxonomic donut charts generated using the MetaPORE bioinformatics analysis pipeline from the Chik1 run (left) and Ebola1 run (right). The total number of reads analyzed is shown in the center of the donut. **(D)** Coverage and pairwise identity plots generated in MetaPORE by mapping reads aligning to CHIKV (left, Chik1 run) or EBOV (right, Ebola1 run) to the closest matching reference genome in the NT database (asterisk). **(E)** Whole-genome phylogeny of CHIKV. Representative CHIKV genome sequences from the Asian-Pacific clade, including the Puerto Rico PR-S6 strain recovered by nanopore and MiSeq sequencing, or all 188 near-complete or complete CHIKV genome sequences available in the NCBI NT database as of March 2015 (inset), are included. Branch lengths are drawn proportionally to the number of nucleotide substitutions per position, and support values are shown for each node. Abbreviations: CHIKV, chikungunya virus; EBOV, Ebola virus; Chik1, chikungunya virus, strain PR-S6 sample; Ebola1, EBOV, strain Lomela-Lokolia16 sample.

476

477 **Figure 3.  Metagenomic identification of HCV from a clinical serum sample**

478 **by nanopore sequencing.  (A)** Timeline of sequencing runs on flowcell #2 with

479 HepC1 sample reloading, plotted as a function of elapsed time in hr since the

480 start of flowcell sequencing.  **(B)** Cumulative number of all sequenced reads

481 (black line) and HCV viral reads (red line), plotted as a function of individual

482 sequencing run time in min.  **(C)** Taxonomic donut charts generated using the

483 MetaPORE bioinformatics analysis pipeline.  The total number of reads analyzed

484 is shown in the center of the donut.  **(D)** Coverage and pairwise identity plots

485 generated in MetaPORE by mapping reads aligning to HCV to the closest

486 matching reference genome in the NT database.  Abbreviations: HCV, hepatitis

487 C virus; HepC1: hepatitis C virus, genotype 1b sample.

488

489 **Figure 4. Metagenomic identification of EBOV from a clinical blood sample**

490 **by nanopore sequencing and MetaPORE real-time bioinformatics analysis.**

491 Nanopore data generated from the Ebola2 library and sequenced on flowcell #3

492 was analyzed in real-time using the MetaPORE bioinformatics analysis pipeline,

493 and compared to corresponding MiSeq data.  **(A)** Timeline of nanopore

494 sequencing runs on flowcell #3 with sample reloading, plotted as a function of

495 elapsed time in hr since the start of flowcell sequencing.  **(B)** Cumulative

496 numbers of all sequenced reads (black line) and target viral reads (red line) from

497 the nanopore run (left) or MiSeq run (right ), plotted as a function of individual

498 sequencing run time in min. **(C)** Taxonomic donut charts generated by real-time

499 MetaPORE analysis of the nanopore reads (left) and post-run analysis of the

500 MiSeq reads (right). The total number of reads analyzed is shown in the center

501 of the donut. Note that only a subset of Illumina reads (n=100,000) were

502 analyzed using the MetaPORE pipeline. **(D)** Coverage and pairwise identity plots

503 generated by MetaPORE analysis of nanopore reads (left) and MiSeq reads

504 (right), generated by mapping reads aligning to EBOV to the closest matching

505 reference genome in the NT database (asterisk). Abbreviations: EBOV, Ebola

506 virus; Ebola2, EBOV, strain Lomela-LokoliaB11 sample.

507

508 **Supplemental Figure 1. MetaPORE analysis of Illumina MiSeq data from**

509 **samples containing CHIKV and EBOV. (A)** Taxonomic donut charts generated

510 from MiSeq Chik1 run data (left) and Ebola1 run data (right) using the

511 MetaPORE bioinformatics analysis pipeline. The total number of MiSeq reads

512 analyzed is shown in the center of the donut. Note that only a subset of reads

513 (n=100,000) were analyzed. **(D)** Coverage and pairwise identity plots generated

514 in MetaPORE by mapping CHIKV reads from the Chik1 run (left) or EBOV reads

515 from the Ebola1 run (right) to the closest matching reference genome in the NT

516 database.

517

518 **Movie. MetaPORE Real-Time Bioinformatics Analysis and Visualization**

519 (clip 0:19 - 10:13, 9.8 min) Detection of EBOV by metagenomic nanopore

520 sequencing and real-time MetaPORE bioinformatics analysis. Raw FAST5 files

521 are uploaded to the Metrichor cloud-based analytics platform for 2D basecalling

522 (right panel). After downloading from Metrichor, basecalled FAST5 reads are

523  collected in batches of 200 reads and automatically processed in real-time by

524  MetaPORE.  Reads corresponding to human host background are

525  computationally subtracted by BLASTn alignment to the human portion of the

526  NCBI NT database.  Following subtraction, remaining reads are then aligned to

527  NCBI NT using BLASTn for taxonomic identification.  The closest matched

528  reference genomes on the basis of hit frequency are shown in donut plots that

529  are updated each minute in real-time (left panel).  Note that the first EBOV read

530  from the Ebola2 sequencing run is detected 3 min 6 sec (3:26) after the start of

531  sequence acquisition (0:19).  (clip 10:21 - 11:51, 1.5 min) Web-based, interactive

532  coverage map and pairwise identity plots, generated in real-time by MetaPORE,

533  enable zooming, highlighting of individual values, outputting of relevant statistical

534  data, and exporting of the graphs in various formats.  The plots shown in the

535  movie correspond to the analyzed data after completion of nanopore sequencing.

536

537  **Data Availability**

538  Nanopore and MiSeq sequencing data, along with sample metadata, have been

539  submitted to NCBI under the following GenBank accession numbers:

540  Sample metadata along with nanopore and MiSeq sequencing data have been

541  submitted to the NCBI for each sample under the following SRA Study

542  accessions: Ebola virus/H.sapiens-wt/COD/2014/Lomela-Lokolia16

543  (SRP057409), Ebola virus/H.sapiens-wt/COD/2014/Lomela-LokoliaB11

544  (SRS933322), Chik1 (SRP057410), and HepC1 (SRP057418).

545

**SOURCE OF FUNDING**

546

547	This study is supported in part by a grant from the National Institutes of

548	Health (R01-HL105704) (CYC) and an UCSF-Abbott Viral Discovery Award

549	(CYC).

550

**COMPETING INTERESTS**

551

552	CYC is the director of the UCSF-Abbott Viral Diagnostics and Discovery

553	Center (VDDC) and receives research support in pathogen discovery from Abbott

554	Laboratories, Inc.  JL and VB are employees of Hologic, Inc.  P. Mbala and P.

555	Mulembakani and BS are employees of Metabiota, Inc.

556

557	AUTHOR'S CONTRIBUTIONS:

558	ALG, and CYC directed the study. ALG performed all nanopore experiments;

559	SNN, GY, SS, and JB provided cDNA libraries and generated Illumina libraries.

560	SF and ALG constructed the MetaPORE pipeline; SF,DS, CYC generated real

561	time visualization; ALG, SF, SNN validated MetaPORE pipeline.  VB, JML, RD,

562	and SS provided chikungunya sample; PM, BSS, JJM provided ebola sample.

563	ALG, SF, SNN, and CYC wrote the manuscript.

564

565

## REFERENCES

1.  Pallen MJ: **Diagnostic metagenomics: potential applications to bacterial, viral and parasitic infections.** *Parasitology* 2014, **141:**1856-1862.

2.  Miller RR, Montoya V, Gardy JL, Patrick DM, Tang P: **Metagenomics for pathogen detection in public health.** *Genome Med* 2013, **5:**81.

3.  Brown JR, Morfopoulou S, Hubb J, Emmett WA, Ip W, Shah D, Brooks T, Paine SM, Anderson G, Virasami A, et al: **Astrovirus VA1/HMO-C: an increasingly recognized neurotropic pathogen in immunocompromised patients.** *Clin Infect Dis* 2015, **60:**881-888.

4.  Naccache SN, Peggs KS, Mattes FM, Phadke R, Garson JA, Grant P, Samayoa E, Federman S, Miller S, Lunn MP, et al: **Diagnosis of neuroinvasive astrovirus infection in an immunocompromised adult with encephalitis by unbiased next-generation sequencing.** *Clin Infect Dis* 2015, **60:**919-923.

5.  Palacios G, Druce J, Du L, Tran T, Birch C, Briese T, Conlan S, Quan PL, Hui J, Marshall J, et al: **A new arenavirus in a cluster of fatal transplant-associated diseases.** *N Engl J Med* 2008, **358:**991-998.

6.  Wilson MR, Naccache SN, Samayoa E, Biagtan M, Bashir H, Yu G, Salamat SM, Somasekar S, Federman S, Miller S, et al: **Actionable diagnosis of neuroleptospirosis by next-generation sequencing.** *N Engl J Med* 2014, **370:**2408-2417.

7.  Briese T, Paweska JT, McMullan LK, Hutchison SK, Street C, Palacios G, Khristova ML, Weyer J, Swanepoel R, Egholm M, et al: **Genetic detection and characterization of Lujo virus, a new hemorrhagic fever-associated arenavirus from southern Africa.** *PLoS Pathog* 2009, **5:**e1000455.

8.  Gire SK, Goba A, Andersen KG, Sealfon RS, Park DJ, Kanneh L, Jalloh S, Momoh M, Fullah M, Dudas G, et al: **Genomic surveillance elucidates Ebola virus origin and transmission during the 2014 outbreak.** *Science* 2014, **345:**1369-1372.

9.  Ashton PM, Nair S, Dallman T, Rubino S, Rabsch W, Mwaigwisya S, Wain J, O'Grady J: **MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island.** *Nat Biotechnol* 2015, **33:**296-300.

10. Goodwin S, Gurtowski J, Ethe-Sayers S, Despande P, Schatz M, McCombie WR: **Oxford Nanopore Sequencing and de novo Assembly of a Eukaryotic Genome.** *bioRxiv* 2015, http://dx.doi.org/10.1101/013490.

11. Kilianski A, Haas JL, Corriveau EJ, Liem AT, Willis KL, Kadavy DR, Rosenzweig CN, Minot SS: **Bacterial and viral identification and differentiation by amplicon sequencing on the MinION nanopore sequencer.** *Gigascience* 2015, **4:**12.

12. Quick J, Ashton P, Calus S, Chatt C, Gossain S, Hawker J, Nair S, Neal K, Nye K, Peters T, et al: **Rapid draft sequencign and real-time**

nanopore sequencing in a hospital outbreak of *Salmonella*. *Genome Biology* 2015, **16**.

13. **The MinION Access Programme - Community** [https://nanoporetech.com/community/the-minion-access-programme]

14. Chiu CY, Bres V, Yu G, Krysztof D, Naccache SN, Lee D, Pfeil J, Linnen JM, Stramer SL: **Emerging genomic assays for identification of chikungunya virus infection in blood donors from Puerto Rico, 2014.** *Emerging Infectious Diseases* 2015, **(in press)**.

15. Naccache SN, Federman S, Veeraraghavan N, Zaharia M, Lee D, Samayoa E, Bouquet J, Greninger AL, Luk KC, Enge B, et al: **A cloud-compatible bioinformatics pipeline for ultrarapid pathogen identification from next-generation sequencing of clinical samples.** *Genome Res* 2014, **24:**1180-1192.

16. Chen EC, Miller SA, DeRisi JL, Chiu CY: **Using a pan-viral microarray assay (Virochip) to screen clinical samples for viral pathogens.** *J Vis Exp* 2011.

17. Greninger AL, Chen EC, Sittler T, Scheinerman A, Roubinian N, Yu G, Kim E, Pillai DR, Guyard C, Mazzulli T, et al: **A metagenomic analysis of pandemic influenza A (2009 H1N1) infection in patients from North America.** *PLoS One* 2010, **5:**e13381.

18. Greninger AL, Naccache SN, Messacar K, Clayton A, Yu G, Somasekar S, Federman S, Stryke D, Anderson C, Yagi S, et al: **A novel outbreak enterovirus D68 strain associated with acute flaccid myelitis cases in the USA (2012-14): a retrospective cohort study.** *Lancet Infect Dis* 2015.

19. Martin M: **Cutadapt removes adapter sequences from high-throughput sequenincg reads.** *EMBnetjournal* 2011, **17:**10-12.

20. Chiu CY: **Viral pathogen discovery.** *Curr Opin Microbiol* 2013, **16:**468-478.

21. Harris RS: **Improved Pairwise Alignment of Genomic DNA.** Pennsylvania State University, 2007.

22. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Genome Project Data Processing S: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25:**2078-2079.

23. Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F, Wilczynski B, de Hoon MJ: **Biopython: freely available Python tools for computational molecular biology and bioinformatics.** *Bioinformatics* 2009, **25:**1422-1423.

24. Needleman SB, Wunsch CD: **A general method applicable to the search for similarities in the amino acid sequence of two proteins.** *J Mol Biol* 1970, **48:**443-453.

25. Maganga GD, Kapetshi J, Berthet N, Kebela Ilunga B, Kabange F, Mbala Kingebeni P, Mondonge V, Muyembe JJ, Bertherat E, Briand S, et al: **Ebola virus disease in the Democratic Republic of Congo.** *N Engl J Med* 2014, **371:**2083-2091.

657    26.  Baize S, Pannetier D, Oestereich L, Rieger T, Koivogui L, Magassouba N,
658         Soropogui B, Sow MS, Keita S, De Clerck H, et al: **Emergence of Zaire**
659         **Ebola virus disease in Guinea.** *N Engl J Med* 2014, **371:**1418-1425.
660    27.  Katoh K, Standley DM: **MAFFT multiple sequence alignment software**
661         **version 7: improvements in performance and usability.** *Mol Biol Evol*
662         2013, **30:**772-780.
663    28.  Huelsenbeck JP, Ronquist F: **MRBAYES: Bayesian inference of**
664         **phylogenetic trees.** *Bioinformatics* 2001, **17:**754-755.
665    29.  Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S,
666         Buxton S, Cooper A, Markowitz S, Duran C, et al: **Geneious Basic: an**
667         **integrated and extendable desktop software platform for the**
668         **organization and analysis of sequence data.** *Bioinformatics* 2012,
669         **28:**1647-1649.
670    30.  Blakemore R, Story E, Helb D, Kop J, Banada P, Owens MR, Chakravorty
671         S, Jones M, Alland D: **Evaluation of the analytical performance of the**
672         **Xpert MTB/RIF assay.** *J Clin Microbiol* 2010, **48:**2495-2501.
673    31.  Fournier PE, Dubourg G, Raoult D: **Clinical detection and**
674         **characterization of bacterial pathogens in the genomics era.** *Genome*
675         *Med* 2014, **6:**114.
676    32.  Kunin V, Copeland A, Lapidus A, Mavromatis K, Hugenholtz P: **A**
677         **bioinformatician's guide to metagenomics.** *Microbiol Mol Biol Rev*
678         2008, **72:**557-578, Table of Contents.
679    33.  Jain M, Fiddes IT, Miga KH, Olsen HE, Paten B, Akeson M: **Improved**
680         **data analysis for the MinION nanopore sequencer.** *Nat Methods* 2015,
681         **12:**351-356.
682    34.  Loman NJ, Watson M: **Successful test launch for nanopore**
683         **sequencing.** *Nat Methods* 2015, **12:**303-304.
684
685

Table 1

| Experiment / Sample | Flow Cell # | Run # | # of Active Pores | Run time (min) | Total reads | Pass reads | Fail reads | Pass / fail rate | Target Virus* |
|---|---|---|---|---|---|---|---|---|---|
| Chik1 | 1 | first run | 345 | 138 | 19,452 | 5,139 | 14,313 | 35.9% | CHIKV |
| Ebola1 | 1 | first run | 105 | 1,022 | 13,090 | 1,831 | 11,259 | 16.3% | EBOV |
| HepC1 | 2 | first run | 171 | 122 | 10,305 | 729 | 9,877 | 7.4% | HCV |
| HepC1 | 2 | reload #1 | 293 | 192 | 26,626 | 2,155 | 25,758 | 8.4% | HCV |
| HepC1 | 2 | reload #2 | 256 | 298 | 32,212 | 1,207 | 31,289 | 3.9% | HCV |
| HepC1 | 2 | reload #3 | 214 | 156 | 14,805 | 287 | 14,275 | 2.0% | HCV |
| Ebola2 | 3 | first run | 397 | 79 | 28,651 | 1,537 | 27,114 | 5.7% | EBOV |
| Ebola2 | 3 | reload #1 | 426 | 222 | 95,861 | 2,899 | 92,962 | 3.1% | EBOV |
| Ebola2 | 3 | reload #2 | 380 | 1,091 | 166,524 | 1,539 | 164,985 | 0.9% | EBOV |
| Ebola2 | 3 | reload #3 | 200 | 1,357 | 44,272 | 34 | 44,238 | 0.1% | EBOV |
| TOTAL | | | | | 451,798 | 17,357 | 436,070 | 4.0% | |

Table S1. Run data for each flow cell.

*abbreviations: CHIKV, chikungunya virus; EBOV, Ebola virus; HCV, hepatitis C virus

**based on average pairwise identity of aligned viral reads to the most closely matched refer

| Avg viral read length | Avg viral read error rate** |
|---|---|
| 375 | 20.6% (8-49%) |
| 317 | 22.0% (12-43%) |
| 497 | |
| 511 | |
| 506 | |
| 479 | 33.1% (24-46%) |
| 366 | |
| 377 | |
| 372 | |
| 326 | 22.3% (8-48%) |
| 391 | 24.3% (8-49%) |

ence sequence

**A**

RNA extraction (+DNase)

1 hr

1 hr

3 hr

3 hr

Reverse transcription and random amplification

**MinION (Oxford Nanopore)**

NGS library preparation

**MiSeq (Illumina)**

NGS library preparation

5 hr

5 hr

Flowcell initialization (20 min)

Library quantitation

6 hr

Nanopore sequencing (~100-300 reads/min)

MiSeq sequencing

>18 hr*

Sequence data acquisition

Metrichor cloud-based basecalling of nanopore reads (~200-300 reads/min)

MetaPORE real-time bioinformatics analysis (~100-200 reads/min)

SURPI bioinformatics analysis (>2 hr)

<6 hr**

TIME TO INITAL DETECTION

>20 hr**

CUMULATIVE TIME

CUMULATIVE TIME

**B**    **MetaPORE Real-Time Bioinformatics Analysis Pipeline**

continuous scanning of nanopore read folder (every 50 reads or 30 sec)

Human BLASTn computational subtraction

NT BLASTn with taxonomic classification

Real-time visualization of taxonomy, coverage maps, and pairwise identity plots

Figure 2

Figure 3

**A**

**FLOWCELL #2 (HCV 1x10^5 cp/mL)**

| HepC1 run start | HepC1 reload #1 | HepC1 reload #2 | HepC1 reload #3 | HepC1 run stop |
|---|---|---|---|---|
| 0 | 2.0 | 5.2 | 10.5 | 13 |

elapsed run time (hr)

**B**



cumulative # of reads vs run minutes

— All reads (n=85,585)
— HCV reads (n=6)

34 min

**C**

HepC1-MinION (All Reads)

85,647

HepC1-MinION (Viruses)

156

*Homo sapiens* (23,652)
Unidentified (61,006)
Viruses (156)
Cutadapt removed (588)
Bacteria (172)
Non-human eukaryote (19)
Other lineage (54)

Hepatitis C virus (6)
Phages (150)

**D**

**Hepatitis C virus subtype 1b (FJ024277, gi|197310699|, 9329 bp)**
Hepatitis C virus subtype 1b isolate HCV-1b/US/BID-V1711/2007, complete genome

1X: 24.1%, 3X: 0.0%

Fold coverage (log)

Supp Figure 1



**A**

Chik1-MiSeq (All Reads) — 100,000

Chik1-MiSeq (Viruses) — 7,626

- *Homo sapiens* (81,947)
- Unidentified (3,390)
- Non-human eukaryote (135)
- Bacteria (13)
- Viruses (7,626)
- Other lineage (15)
- Cutadapt removed (6,874)

- Chikungunya virus (7,626)

**B**

Ebola1-MiSeq (All Reads) — 100,000

Ebola1-MiSeq (Viruses) — 840

- *Homo sapiens* (56,504)
- Unidentified (2,616)
- Non-human eukaryote (288)
- Bacteria (1,189)
- Viruses (840)
- Other lineage (722)
- Cutadapt removed (37,841)

- Zaire ebolavirus (763)
- Phages (77)

**C**

Chikungunya virus (gi|615794507|, 12011 bp)

*Chikungunya virus strain 99659, complete genome*

1X: 99.9%, 3X: 99.9%

99.5%

**D**

Zaire ebolavirus (gi|733962878|, 18941 bp)

*Zaire ebolavirus isolate Ebola virus/H.sapiens-wt/COD/2014/Lomela-Lokolia16, complete genome*

1X: 99.7%, 3X: 98.9%

99.2%