| | |
|---|---|
| Article Type: | Invited Methods Article - JoVE Produced Video |
| Manuscript Number: | JoVE60906R2 |
| Full Title: | Validating whole genome Nanopore sequencing, using Usutu virus as an example |
| Section/Category: | JoVE Genetics |
| Keywords: | Nanopore sequencing;  R10 flowcell;  USUV;  arboviruses;  whole genome sequencing |
| Corresponding Author: | Bas Oude Oude Munnink<br>Erasmus MC<br>Rotterdam, Zuid-Holland NETHERLANDS |
| Corresponding Author's Institution: | Erasmus MC |
| Corresponding Author E-Mail: | b.oudemunnink@erasmusmc.nl |
| Order of Authors: | Bas B. Oude Munnink |
| | David F. Nieuwenhuijse |
| | Reina S. Sikkema |
| | Marion Koopmans |
| Additional Information: | |
| Question | Response |
| Please indicate whether this article will be Standard Access or Open Access. | Open Access (US$4,200) |
| Please indicate the **city, state/province, and country** where this article will be **filmed**. Please do not use abbreviations. | Rotterdam, Zuid-Holland, The Netherlands |

1  **TITLE:**
2  Validating Whole Genome Nanopore Sequencing, using Usutu Virus as an Example
3
4  **AUTHORS AND AFFILIATIONS:**
5  Bas B. Oude Munnink[1], David F. Nieuwenhuijse[1], Reina S. Sikkema[1] and Marion Koopmans[1]
6
7  [1]ErasmusMC, Department of Viroscience, WHO Collaborating Centre for Arbovirus and Viral
8  Hemorrhagic Fever Reference and Research, Rotterdam, The Netherlands
9
10  **Email Addresses of Co-Authors:**
11  David F. Nieuwenhuijse          (d.nieuwenhuijsse@erasmusmc.nl)
12  Reina S. Sikkema               (r.sikkema@erasmusmc.nl)
13  Marion Koopmans                (m.koopmans@erasmusmc.nl)
14
15  **Corresponding Author:**
16  Bas B. Oude Munnink             (b.oudemunnink@erasmusmc.nl)
17
18  **KEYWORDS:**
19  Nanopore, sequencing, R10 flowcell, USUV, arboviruses, whole genome sequencing
20
21  **SUMMARY:**
22  We previously validated a protocol for amplicon-based whole genome Usutu virus (USUV)
23  sequencing on a nanopore sequencing platform. Here, we describe the methods used in more
24  detail and determine the error rate of the nanopore R10 flow cell.
25
26  **ABSTRACT:**
27  Whole genome sequencing can be used to characterize and to trace viral outbreaks. Nanopore-
28  based whole genome sequencing protocols have been described for several different viruses.
29  These approaches utilize an overlapping amplicon-based approach which can be used to target a
30  specific virus or group of genetically related viruses. In addition to confirmation of the virus
31  presence, sequencing can be used for genomic epidemiology studies, to track viruses and unravel
32  origins, reservoirs and modes of transmission. For such applications, it is crucial to understand
33  possible effects of the error rate associated with the platform used. Routine application in clinical
34  and public health settings require that this is documented with every important change in the
35  protocol. Previously, a protocol for whole genome Usutu virus sequencing on the nanopore
36  sequencing platform was validated (R9.4 flowcell) by direct comparison to Illumina sequencing.
37  Here, we describe the method used to determine the required read coverage, using the
38  comparison between the R10 flow cell and Illumina sequencing as an example.
39
40  **INTRODUCTION:**
41  Fast developments in third generation sequence technologies allows us to move forward towards
42  close to real-time sequencing during viral outbreaks. This timely availability of genetic
43  information can be useful to determine the origin and evolution of viral pathogens. Gold
44  standards in the fields of next generation sequencing however, are still the second-generation

45 sequencers. These techniques rely on specific and time-consuming techniques like clonal
46 amplification during an emulsion PCR or clonal bridge amplification. The third-generation
47 sequencers are cheaper, hand-held and come with simplified library preparation methodologies.
48 Especially the small size of the sequence device and the low purchase price makes it an
49 interesting candidate for deployable, fieldable sequencing. This could for instance be seen during
50 the Ebola virus outbreak in Sierra Leone and during the ongoing arbovirus outbreak investigations
51 in Brazil[1–3]. However, the reported high error rate[4] might limit the applications for which
52 nanopore sequencing can be used.

54 Nanopore sequencing is evolving quickly. New products are available in the market on a regular
55 basis. Examples of this are for instance the 1D squared kits which enables sequencing of both
56 strands of the DNA molecule, thereby boosting the accuracy of the called bases[5] and the
57 development of the R10 flow cell which measures the change in current at two different instances
58 in the pore[6]. In addition, improved bio-informatic tools like improvements in basecalling will
59 improve the accuracy of basecalling[7]. One of the most frequently used basecallers, (e.g.,
60 Albacore), has been updated at least 12 times in a 9-month time period[5]. Recently, the
61 manufacturer also released a novel basecaller called flip-flop, which is implemented in the
62 default nanopore software[8]. Together, all of these improvements will lead to more accurate
63 sequences and will decrease the error rate of the nanopore sequencer.

65 Usutu virus (USUV) is a mosquito-borne arbovirus of the family *Flaviviridae* and it has a positive-
66 stranded RNA genome of around 11,000 nucleotides. USUV mainly affects great grey owls and
67 blackbirds[9, 10], although other bird species are also susceptible to USUV infection[11]. Recently,
68 USUV was also identified in rodents and shrews although their potential role in transmission of
69 the virus remains unknown[12]. In humans, asymptomatic infections have been described in blood
70 donors[13–16] while USUV infections also have been reported to be associated with encephalitis or
71 meningo-encephalitis[17, 18]. In the Netherlands, USUV was first detected in wild birds in 2016[10]
72 and in asymptomatic blood donors in 2018[14]. Since the initial detection of USUV, outbreaks have
73 been reported during the subsequent years and surveillance, including whole genome
74 sequencing, is currently ongoing to monitor the emerge and spread of an arbovirus in a previously
75 naïve population.

77 Similar to what has been described for other viruses, such as Ebola virus, Zika virus and yellow
78 fever virus[3, 19, 20], we have developed a primer set to sequence full length USUV[21]. This
79 polymerase chain reaction (PCR)-based approach allows for the recovery of full length USUV
80 genomes from highly host-contaminated sample types like brain samples in samples up to a Ct
81 value of around 32. Benefits of an amplicon-based sequencing approach are a higher sensitivity
82 compared to metagenomic sequencing and a higher specificity. Limitations of using an amplicon-
83 based approach are that the sequences should be similar in order to design primers fitting all
84 strains and that primers are designed on our current knowledge about the virus diversity.

86 Given the constant developments and improvements in third generation sequencing, there is a
87 need to evaluate the error rate of the sequencer on a regular basis. Here, we describe a method
88 to evaluate the performance of nanopore directly against Illumina sequencing using USUV as an

89  example. This method is applied to sequences generated with the latest R10 flow cell and
90  basecalling is performed with the latest version of the flip-flop basecaller.
91
92  **PROTOCOL:**
93
94  NOTE: List of software tools to be used: usearch v11.0.667; muscle v3.8.1551; porechop 0.2.4;
95  cutadapt 2.5; minimap2 2.16-r922; samtools 1.9; trimmomatic 0.39; bbmap 38.33; spades
96  v3.13.1; kma-1.2.8
97
98  **1. Primer design**
99
100 1.1. Start with downloading or retrieving a set of relevant reference whole genome sequences
101 from public or private data collections. For instance, retrieve all full length USUV genomes
102 (taxid64286) from the NCBI database[22]. USUV encodes a genome of around 11,000 nucleotides
103 so only retrieve the sequences with a sequence length of 8,000–12,000 nucleotides. Do this using
104 the following search entry:
105 -     *taxid64286[Organism:noexp] AND 8000[SLEN]:12000[SLEN].*
106
107 1.1.1. Click on **Send to | Complete Record | File**; use Format = FASTA and create the File.
108
109 1.2. To downsize the set of reference sequences, remove duplicate sequences or sequences with
110 over 99% nucleotide identity from the dataset. Do this using the cluster fast option from
111 usearch[23]. On the command line enter:
112 -     *usearch -cluster_fast All_USUV.fasta -id 0.99 -centroids All_USUV_dedup.fasta*
113
114 1.3. To generate the primers, sequences need to be aligned. This is done using MUSCLE[24]. On the
115 command line enter:
116 -     *muscle -in All_USUV_dedup.fasta -out All_USUV_dedup_aligned.fasta -log log_muscle.txt*
117
118 NOTE: It is essential to manually inspect the alignment to check for discrepancies. These can be
119 manually corrected if needed and the ends can be trimmed according to the length of most whole
120 genome sequences.
121
122 1.4. Primal is used to make a draft selection of the primers which can be used for full length
123 amplicon        sequencing[19].      Upload       the       alignment       to       the       primal       website
124 (http://primal.zibraproject.org/) and select the preferred amplicon length and overlap length
125 between the different amplicons. Go to primal.zibraproject.org, fill in the **Scheme name**, upload
126 the aligned fasta file, select the amplicon length, overlap size, and generate the scheme.
127
128 1.5. Align the complete set of available complete USUV sequences (not the downsized or
129 deduplicated set). On the command line enter:
130 -     *muscle -in All_USUV.fasta -out All_USUV_aligned.fasta -log log_muscle.txt*
131
132 NOTE: Map the generated primers against the complete alignment (do not use the deduplicated

133 alignment), manually correct errors and include a maximum of 5 degenerative primer positions.
134
135 **2. Multiplex PCR**
136
137 2.1. Perform the multiplex PCR using the designed primers and nanopore and Illumina
138 sequencing. The multiplex PCR for USUV was performed as previous described[19, 21].
139
140 2.2. Perform basecalling with flip-flop version 3.0.6.6+9999d81.
141
142 **3. Data analysis to generate consensus sequences from nanopore data**
143
144 3.1. Several samples can be multiplexed on a single nanopore sequencing run. After performing
145 the sequence run, demultiplex the nanopore data. Use Porechop[25] for this. To prevent
146 contamination and enhance accuracy, use the *require_two_barcodes* flag. On the command line
147 enter:
148 - *porechop -i Run_USUV.fastq -o Run_USUV_demultiplex --require_two_barcodes*
149
150 3.2. After demultiplexing, remove primer sequences (indicated in the file Primers_Usutu.fasta in
151 both orientations) using cutadapt[26]. In addition, remove sequences with a length shorter than 75
152 nucleotides. The primers have to be removed since they can introduce artificial biases in the
153 consensus sequence. On the command line enter:
154 - *cutadapt -b file:Primers_USUV.fasta -o BC01_trimmed.fastq BC01.fastq -m 75*
155
156 3.3. Demultiplexed sequence reads can be mapped against a panel of distinct reference strains
157 using minimap2[27] and a consensus sequence can be generated using samtools[28]. Follow the
158 example below which shows the procedure of a reference-based alignment and the consensus
159 sequence generation of one sample: BC01. On the command line enter:
160 - *minimap2 -ax map-ont Random_Refs_USUV.fasta BC01_trimmed.fastq > BC01.bam*
161 - *samtools sort BC01.bam > BC01_sorted.bam*
162 - *bcftools mpileup -Ou -f Random_Refs_USUV.fasta BC01_sorted.bam | bcftools call -mv -*
163 *Oz -o BC01.vcf.gz*
164 - *bcftools index BC01.vcf.gz*
165 - *cat Random_Refs_USUV.fasta | bcftools consensus BC01.vcf.gz > BC01_consensus.fasta*
166
167 3.4. For reference-based alignments it is essential that a closely related reference sequence is
168 used. Therefore, perform a BlastN search with the generated consensus sequence to identify the
169 closest reference strain. After that, repeat the reference-based alignment with the closest
170 reference strain as reference (step 3.3 and 3.4). On the command line enter:
171 - *minimap2 -ax map-ont Ref_USUV_BC01.fasta BC01_trimmed.fastq > BC01_ref.bam*
172 - *samtools sort BC01_ref.bam > BC01_sorted_ref.bam*
173 - *bcftools mpileup -Ou -f Ref_USUV_BC01.fasta BC01_sorted_ref.bam | bcftools call -mv -*
174 *Oz -o BC01_ref.vcf.gz*
175 - *bcftools index BC01_ref.vcf.gz*
176 - *cat Ref_USUV_BC01.fasta | bcftools consensus BC01_ref.vcf.gz >*

177    *BC01_ref_consensus.fasta*

178

179    **4. Analysis of the Illumina data**

180

181    4.1. These sequences are automatically demultiplexed after sequencing. Reads can be quality
182    controlled using trimmomatic[29]. For paired-end Illumina sequences, use the commonly used cut-
183    off median PHRED score of 33 and a minimal read length of 75 to get accurate, high quality reads.
184    On the command line enter:
185    -      *trimmomatic  PE  -phred33  9_S9_L001_R1_001.fastq.gz  9_S9_L001_R2_001.fastq.gz*
186    *9_1P.fastq  9_1U.fastq  9_2P.fastq  9_2U.fastq  LEADING:3  TRAILING:3  SLIDINGWINDOW:3:15*
187    *MINLEN:75*

188

189    4.2. Remove primers (indicated in the file Primers_Usutu.fasta in both orientations), since they
190    can introduce artificial biases, using cutadapt[26]. In addition, remove sequences with a length
191    shorter than 75 nucleotides using the commands below. On the command line enter:
192    -      *cutadapt  -b  file:Primers_USUV.fasta  -o  9_1P_trimmed.fastq  -p  9_2P_trimmed.fastq*
193    *9_1P.fastq 9_2P.fastq -m 75*

194

195    4.3. Before de novo assembly, the sequence reads can be normalized for an even coverage across
196    the genome. This is essential since de novo assemblers like SPAdes take the read coverage into
197    account when assembling sequence reads. Normalize reads to a read coverage of 50 using
198    BBNorm from the BBMap package[30]. On the command line enter:
199    -      *bbmap/bbnorm.sh    target=50    in=9_1P_trimmed.fastq    in2=9_2P_trimmed.fastq*
200    *out=Sample9_FW_norm.fastq out2= Sample9_RE_norm.fastq*

201

202    4.4. The normalized reads are de novo assembled using SPAdes[31]. Default settings are used for
203    the assembly using all different kmers (21, 33, 55, 77, 99 and 127). On the command line enter:
204    -      *spades.py -k 21,33,55,77,99,127 -o Sample9 -1 Sample9.qc.f.fq -2 Sample9.qc.r.fq*

205

206    4.5. Map the QC reads against the obtained consensus sequence using minimap2 and programs
207    like Geneious, Bioedit or Ugene to curate the alignment. It is important to check the beginning
208    and the end of the contig.

209

210    4.5.1. Align the QC reads against the obtained consensus sequencing using minimap2.

211

212    4.5.2. Import the alignment in Geneious/Bioedit/UGene.

213

214    4.5.3. Manually inspect, correct and curate especially the beginning and the end of the genome.

215

216    **5. Determining the required read coverage to compensate for the error profile in nanopore**
217    **sequencing using Illumina data as gold standard**

218

219    5.1. Select sequence reads mapping to one amplicon, in this case amplicon 26. Subsequently,
220    map the nanopore reads against this amplicon using minimap2. Use Samtools to select only the

221 reads mapping to amplicon 26 and to convert the bam file into fastq. On the command line enter:
222 -      *minimap2 -ax map-ont -m 150 Amplicon26.fasta BC01_trimmed.fastq > BC01.bam*
223 -      *samtools view -b -F 4 BC01.bam > BC01_mapped.bam*
224 -      *samtools bam2fq BC01_mapped.bam | seqtk seq - -> BC01_mapped.fastq*
225

226 5.2. Randomly select subsets of for instance 200 sequence reads one thousand times. For
227 example, changing it to 10 will result in the random selection of one thousand times a subset of
228 10 sequence reads. The script is provided as **Supplementary File 1**. On the command line enter:
229 -      *python Random_selection.py*
230

231 5.3. All randomly selected sequence reads are aligned to amplicon 26. Use KMA[32] to map the
232 sequence reads and to immediately generate a consensus sequence. Use optimized settings for
233 nanopore sequencing, indicated by the -bcNano flag. On the command line enter:
234 *kma index -i Amplicon26.fasta*
235 *for file in random_sample*; do*
236      *sampleID=${file%.fastq}*
237      *kma -i ${sampleID}.fastq -o ${sampleID} -t_db Amplicon26.fasta –mem_mode -mp 5 -mrs*
238 *0.0 -bcNano*
239 *done*
240

241 5.4. Inspect the generated consensus sequences on the command line using:
242 -      *cat *.fsa > All_genomes.fsa*
243 -      *minimap2 -ax map-ont Amplicon26.fasta All_genomes.fsa > All_genomes.bam*
244 -      *samtools sort All_genomes.bam > All_genomes_sorted.bam*
245 -      *samtools stats All_genomes_sorted.bam > stats.txt*
246

247 5.4.1. The error rate is displayed in the stats.txt under the heading **error rate #mismatches /**
248 **bases mapped**. Display it on the screen with the following command:
249 -      *grep ^SN stats.txt | cut -f 2-*
250

251 5.4.2. The amount of indels is displayed under the heading **#Indels per cycle**. Display it on the
252 screen with the following command:
253 -      *grep ^IC stats.txt | cut -f 2-*
254

255 **REPRESENTATIVE RESULTS:**
256 Recently, a new version of the flow cell version (R10) was released and offered improvements to
257 the basecaller used to convert the electronic current signal to DNA sequences (so-called flip-flop
258 basecaller). Therefore, we have re-sequenced USUV from brain tissue of an USUV-positive owl
259 which was previously sequenced on a R9.4 flow cell and on an Illumina Miseq instrument[21]. Here,
260 we described the method used to determine the required read coverage for reliable consensus
261 calling by direct comparison to Illumina sequencing.
262

263 Using the newer flow cell in combination with the basecaller flip-flop we show that a read
264 coverage of 40x results in identical results as compared to Illumina sequencing. A read coverage

265 of 30x results in an error rate of 0.0002% which corresponds to one error in every 585,000
266 nucleotides sequenced, while a read coverage of 20x results in one error in every 63,529
267 nucleotides sequenced. A read coverage of 10x results in one error in every 3,312 nucleotides
268 sequenced, meaning that over three nucleotides per full USUV genome are being called wrong.
269 With a read coverage above 30x, no indels were observed. A read coverage of 20x resulted in the
270 detection of one indel position while a read coverage of 10x resulted in indels in 29 positions. An
271 overview of the error rate using different read coverage cut-offs is shown in **Table 1**.
272
273 **FIGURE AND TABLE LEGENDS:**
274
275 **Table 1. Overview of the error rate of nanopore sequencing.** Each iteration represents one
276 thousand random samples.
277
278 **DISCUSSION:**
279 Nanopore sequencing is constantly evolving and therefore there is a need for methods to monitor
280 the error rate. Here, we describe a workflow to monitor the error rate of the nanopore
281 sequencer. This can be useful after the release of a new flow cell, or if new releases of the
282 basecalling are released. However, this can also be useful for users who want to set-up and
283 validate their own sequencing protocol.
284
285 Different software and alignment tools can yield different results[33]. In this manuscript, we aimed
286 to use freely available software packages which are commonly used, and which have clear
287 documentation. In some cases, preference might be given to commercial tools, which generally
288 have a more user-friendly interfaces but have to be paid for. In the future, this method can be
289 applied to the same sample in case big modifications in sequence technology or basecalling
290 software are introduced Preferentially this should be done after each update of the basecaller or
291 flowcell, however given the speed of the current developments this can be also been done only
292 after major updates.
293
294 The reduction in the error rate in sequencing allows for a higher number of samples to be
295 multiplexed. Thereby, nanopore sequencing is getting closer to replacing conventional real time
296 PCRs for diagnostic assays, which is already the case for influenza virus diagnostics. In addition,
297 the reduction of the error rate increases the usability of this technique sequencing, for instance
298 for the determination of minor variants and for high-throughput unbiased metagenomic
299 sequencing.
300
301 A critical step in the protocol is that close, reliable reference sequences need to be available. The
302 primers are based on the current knowledge about virus diversity and might need to be updated
303 every once in a while. Another critical point when setting up an amplicon-based sequencing
304 approach is the balancing of the primer concentration to get an even balance in amplicon depth.
305 This enables the multiplexing of more samples on a sequence run and results in a significant cost
306 reduction.
307
308 **ACKNOWLEDGMENTS:**

312  **DISCLOSURES:**
313  The authors have nothing to disclose.
314

315  **REFERENCES:**
316  1.     Faria, N. R. et al. Establishment and cryptic transmission of Zika virus in Brazil and the
317  Americas. *Nature*. **546** (7658), 406–410 (2017).
318  2.     Bonaldo, M. C. et al. Genome analysis of yellow fever virus of the ongoing outbreak in
319  Brazil reveals polymorphisms. *Memórias do Instituto Oswaldo Cruz*. **112** (6), 447–451 (2017).
320  3.     Faria, N. R. et al. Genomic and epidemiological monitoring of yellow fever virus
321  transmission potential. *bioRxiv*. 299842 (2018).
322  4.     Magi, A., Giusti, B., Tattini, L. Characterization of MinION nanopore data for resequencing
323  analyses. *Briefings in Bioinformatics*. **18** (6), bbw077 (2016).
324  5.     Rang, F. J., Kloosterman, W.P., de Ridder, J. From squiggle to basepair: computational
325  approaches for improving nanopore sequencing read accuracy. *Genome Biology*. **19** (1), 90
326  (2018).
327  6.     Nanopore Store, R10 flow cells.
328  7.     Wick, R. R., Judd, L. M., Holt, K.E. Performance of neural network basecalling tools for
329  Oxford Nanopore sequencing. *Genome Biology*. **20** (1), 129 (2019).
330  8.     GitHub - nanoporetech/flappie: Flip-flop basecaller for Oxford Nanopore reads.
331  9.     Lühken, R. et al. Distribution of Usutu Virus in Germany and Its Effect on Breeding Bird
332  Populations. *Emerging Infectious Diseases*. **23** (12), 1994–2001 (2017).
333  10.    Cadar, D. et al. Widespread activity of multiple lineages of Usutu virus, Western Europe,
334  2016. *Eurosurveillance*. **22** (4) (2017).
335  11.    Becker, N. et al. Epizootic emergence of Usutu virus in wild and captive birds in Germany.
336  *PLoS ONE*. **7** (2) (2012).
337  12.    Diagne, M. et al. Usutu Virus Isolated from Rodents in Senegal. *Viruses*. **11** (2), 181 (2019).
338  13.    Bakonyi, T. et al. Usutu virus infections among blood donors, Austria, July and August
339  2017 – Raising awareness for diagnostic challenges. *Eurosurveillance*. **22** (41) (2017).
340  14.    Zaaijer, H. L., Slot, E., Molier, M., Reusken, C.B.E.M., Koppelman, M.H.G.M. Usutu virus
341  infection in Dutch blood donors. *Transfusion*. trf.15444 (2019).
342  15.    Cadar, D. et al. Blood donor screening for West Nile virus (WNV) revealed acute Usutu
343  virus (USUV) infection, Germany, September 2016. *Eurosurveillance*. **22** (14), 30501 (2017).
344  16.    Pierro, A. et al. Detection of specific antibodies against West Nile and Usutu viruses in
345  healthy blood donors in northern Italy, 2010–2011. *Clinical Microbiology and Infection*. **19** (10),
346  E451–E453 (2013).
347  17.    Pecorari, M. et al. First human case of Usutu virus neuroinvasive infection, Italy, August-
348  September 2009. *Euro surveillance : bulletin européen sur les maladies transmissibles = European
349  Communicable Disease Bulletin*. **14** (50) (2009).
350  18.    Simonin, Y. et al. Human Usutu Virus Infection with Atypical Neurologic Presentation,
351  Montpellier, France, 2016. *Emerging Infectious Diseases*. **24** (5), 875–878 (2018).
352  19.    Quick, J. et al. Multiplex PCR method for MinION and Illumina sequencing of Zika and

353  other virus genomes directly from clinical samples. *Nature Protocols*. **12** (6), 1261–1276 (2017).
354  20.      Quick, J. et al. Real-time, portable genome sequencing for Ebola surveillance. *Nature*. **530**
355  (7589), 228–232 (2016).
356  21.      Oude Munnink, B. B. et al. Towards high quality real-time whole genome sequencing
357  during outbreaks using Usutu virus as example. *Infection, Genetics and Evolution*. **73**, 49–54
358  (2019).
359  22.      Benson, D. A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J., Sayers, E.W. GenBank. *Nucleic*
360  *Acids Research*. **38** (Database issue), D46-51 (2010).
361  23.      Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*.
362  **26** (19), 2460–2461 (2010).
363  24.      Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
364  throughput. *Nucleic Acids Research*. **32** (5), 1792–1797 (2004).
365  25.      R. R. Wick GitHub - rrwick/Porechop: adapter trimmer for Oxford Nanopore reads.
366  26.      Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing
367  reads. *EMBnet.journal*. **17** (1), 10 (2011).
368  27.      Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. **34** (18),
369  3094–3100 (2018).
370  28.      Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. **25** (16),
371  2078–2079 (2009).
372  29.      Bolger, A. M., Lohse, M., Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
373  data. *Bioinformatics (Oxford, England)*. **30** (15), 2114–20 (2014).
374  30.      BBMap download | SourceForge.net.
375  31.      Bankevich, A. et al. SPAdes: a new genome assembly algorithm and its applications to
376  single-cell sequencing. *J Comput Biol*. **19** (5), 455–477 (2012).
377  32.      Clausen, P. T. L. C., Aarestrup, F.M., Lund, O. Rapid and precise alignment of raw reads
378  against redundant databases with KMA. *BMC Bioinformatics*. **19** (1), 307 (2018).
379  33.      Brinkmann, A. et al. Proficiency Testing of Virus Diagnostics Based on Bioinformatics
380  Analysis of Simulated In Silico High-Throughput Sequencing Data Sets. *Journal of Clinical*
381  *Microbiology*. **57** (8) (2019).
382

Figure

| Coverage | Errors iteration 1 | Error rate iteration 1 | Indels: | Errors iteration 2 | Error rate iteration 2 | Indels: |
|---|---|---|---|---|---|---|
| 10× | 100 | 0.0274% | 4 | 116 | 0.0297% | 18 |
| 20× | 4 | 0.0010% | 0 | 6 | 0.0015% | 1 |
| 30x | 2 | 0.0005% | 0 | 0 | 0.0000% | 0 |
| 40x | 0 | 0.0000% | 0 | 0 | 0.0000% | 0 |
| 50× | 0 | 0.0000% | 0 | 0 | 0.0000% | 0 |

| Coverage | Errors iteration 1 | Error rate iteration 1 | Indels: | Errors iteration 2 | Error rate iteration 2 | Indels: |
|---|---|---|---|---|---|---|

| Errors iteration 3 | Error rate iteration 3 | Indels: |
|---|---|---|
| 110 | 0.0282% | 7 |
| 7 | 0.0018% | 0 |
| 0 | 0.0000% | 0 |
| 0 | 0.0000% | 0 |
| 0 | 0.0000% | 0 |

Table

| Coverage | Errors iteration 1 | Error rate iteration 1 | Indels: | Errors iteration 2 | Error rate iteration 2 | Indels: |
|---|---|---|---|---|---|---|
| 10× | 100 | 0.0274% | 4 | 116 | 0.0297% | 18 |
| 20× | 4 | 0.0010% | 0 | 6 | 0.0015% | 1 |
| 30x | 2 | 0.0005% | 0 | 0 | 0.0000% | 0 |
| 40x | 0 | 0.0000% | 0 | 0 | 0.0000% | 0 |
| 50× | 0 | 0.0000% | 0 | 0 | 0.0000% | 0 |

| Coverage | Errors iteration 1 | Error rate iteration 1 | Indels: | Errors iteration 2 | Error rate iteration 2 | Indels: |
|---|---|---|---|---|---|---|

| Errors iteration 3 | Error rate iteration 3 | Indels: |
|---|---|---|
| 110 | 0.0282% | 7 |
| 7 | 0.0018% | 0 |
| 0 | 0.0000% | 0 |
| 0 | 0.0000% | 0 |
| 0 | 0.0000% | 0 |

| Errors iteration 3 | Error rate iteration 3 | Indels: |
|---|---|---|

| Name of Material/ Equipment | Company | Catalog Number | Comments/Description |
|---|---|---|---|
| Agencourt AMPure XP beads | Beckman Coulter | A63881 | |
| dNTPs | Qiagen | 201900 | |
| FLO-MIN106 R10 flowcell | Nanopore | R10 flowcell | |
| KAPA Hyperplus libarary preparation kit | Roche | 7962436001 | |
| Library Loading Bead Kit | Nanopore | EXP-LLB001 | |
| Ligation Sequencing Kit 1D | Nanopore | SQK-LSK109 | |
| Native Barcoding Kit 1D 1-12 | Nanopore | EXP-NBD103 | |
| Native Barcoding Kit 1D 13-24 | Nanopore | EXP-NBD104 | |
| NEB Blunt/TA Ligase Master Mix | NEB | M0367S | |
| NEB Next Quick Ligation Module | NEB | E6056 | |
| NEB Next Ultra II End Repair / dA-Tailing Module | NEB | E7546S | |
| Protoscript II Reverse Transcriptase | NEB | M0368X | |
| Q5 High-Fidelity polymerase | NEB | M0491 | |
| Qubit dsDNA HS Assay kit | Thermo Fisher | Q32851 | |
| Random Primers | Promega | C1181 | |
| RNAsin Ribonuclease Inhibitor | Promega | N2111 | |

Dear editor,

I have addressed the comments in the resubmitted file. However, the commercial language cannot completely be removed. The manuscript describes the comparison of two different sequence platform which has to be called by name in order to make it understandable for the reader.

Best regards,

Bas Oude Munnink

こ

```
(base) bas:Test_JOVE bom86$ porechop -i Run_USUV.fastq -o Run_USUV_demultiplex --require_two_barcodes

Loading reads
Run_USUV.fastq
198,081 reads loaded


Looking for known adapter sets
550 / 10,000 (5.5%)
```

```
(base) bas:Test_JOVE bom86$ cutadapt -b file:Primers_USUV.fasta -o BC01_trimmed.fastq BC01.fastq -m 75
This is cutadapt 2.5 with Python 3.7.3
Command line parameters: -b file:Primers_USUV.fasta -o BC01_trimmed.fastq BC01.fastq -m 75
Processing reads on 1 core in single-end mode ...
```

```
(base) bas:JOVE bom86$ minimap2 -ax map-ont Random_Refs_USUV.fasta BC01_trimmed.fastq > BC01.bam
[M::mm_idx_gen::0.004*1.70] collected minimizers
[M::mm_idx_gen::0.007*2.22] sorted minimizers
[M::main::0.007*2.22] loaded/built the index for 7 target sequence(s)
[M::mm_mapopt_update::0.008*2.15] mid_occ = 11
[M::mm_idx_stat] kmer size: 15; skip: 10; is_hpc: 0; #seq: 7
[M::mm_idx_stat::0.008*2.11] distinct minimizers: 6651 (62.25% are singletons); average occurrences: 2.151; average spacing: 5.368
```

```
(base) bas:Test_JOVE bom86$ minimap2 -ax map-ont Ref_USUV_BC01.fasta BC01.qc.fq > BC01_ref.bam
[M::mm_idx_gen::0.001*2.82] collected minimizers
[M::mm_idx_gen::0.003*2.85] sorted minimizers
[M::main::0.003*2.82] loaded/built the index for 1 target sequence(s)
[M::mm_mapopt_update::0.003*2.74] mid_occ = 3
[M::mm_idx_stat] kmer size: 15; skip: 10; is_hpc: 0; #seq: 1
[M::mm_idx_stat::0.003*2.66] distinct minimizers: 2047 (99.85% are singletons); average occurrences: 1.001; average spacing: 5.398
ERROR: failed to open file 'BC01.qc.fq'
[M::main] Version: 2.16-r922
[M::main] CMD: minimap2 -ax map-ont Ref_USUV_BC01.fasta BC01.qc.fq
[M::main] Real time: 0.004 sec; CPU: 0.009 sec; Peak RSS: 0.002 GB
(base) bas:Test_JOVE bom86$ samtools sort BC01_ref.bam > BC01_sorted_ref.bam
(base) bas:Test_JOVE bom86$ bcftools mpileup -Ou -f Ref_USUV_BC01.fasta BC01_sorted_ref.bam | bcftools call -mv -Oz -o BC01_ref.vcf.gz
Note: none of --samples-file, --ploidy or --ploidy-file given, assuming all sites are diploid
(mpileup) 1 samples in 1 input files
(base) bas:Test_JOVE bom86$ bcftools index BC01_ref.vcf.gz
(base) bas:Test_JOVE bom86$ cat Ref_USUV_BC01.fasta | bcftools consensus BC01_ref.vcf.gz > BC01_ref_consensus.fasta
Note: the --sample option not given, applying all records regardless of the genotype
(base) bas:Test_JOVE bom86$
(base) bas:Test_JOVE bom86$
```

```
(base) bas:JOVE bom86$ minimap2 -ax map-ont -m 150 Amplicon26.fasta BC01_trimmed.fastq > BC01.bam
[M::mm_idx_gen::0.001*4.68] collected minimizers
[M::mm_idx_gen::0.001*3.72] sorted minimizers
[M::main::0.001*3.63] loaded/built the index for 1 target sequence(s)
[M::mm_mapopt_update::0.001*3.54] mid_occ = 2
[M::mm_idx_stat] kmer size: 15; skip: 10; is_hpc: 0; #seq: 1
[M::mm_idx_stat::0.001*3.46] distinct minimizers: 74 (100.00% are singletons); average occurrences: 1.000; average spacing: 5.270
[M::worker_pipeline::4.448*1.81] mapped 531466 sequences
[M::main] Version: 2.16-r922
[M::main] CMD: minimap2 -ax map-ont -m 150 Amplicon26.fasta BC01_trimmed.fastq
[M::main] Real time: 4.450 sec; CPU: 8.060 sec; Peak RSS: 0.562 GB
(base) bas:JOVE bom86$ samtools view -b -F 4 BC01.bam > BC01_mapped.bam
(base) bas:JOVE bom86$ samtools bam2fq BC01_mapped.bam | seqtk seq - -> BC01_mapped.fastq
[M::bam2fq_mainloop] discarded 0 singletons
[M::bam2fq_mainloop] processed 1834 reads
(base) bas:JOVE bom86$ ▊
```

```
(base) bas:JOVE bom86$ python Random_selection.py
(base) bas:JOVE bom86$ █
```

```
(base) bas:JOVE bom86$ for file in random_sample*; do
> sampleID=${file%.fastq}
> kma -i ${sampleID}.fastq -o ${sampleID} -t_db Amplicon26.fasta -mem_mode -mp 5 -mrs 0.0 -bcNano
> done
```
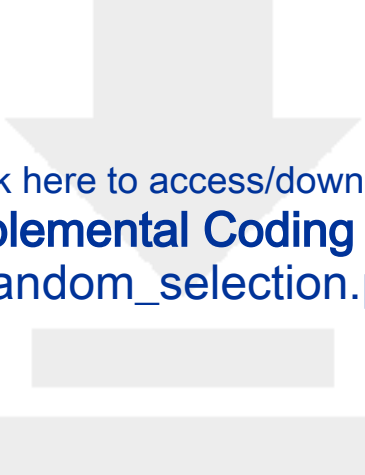
```
(base) bas:JOVE bom86$ cat *.fsa > All_genomes.fsa
(base) bas:JOVE bom86$ minimap2 -ax map-ont Amplicon26.fasta All_genomes.fsa > All_genomes.bam
[M::mm_idx_gen::0.001*2.41] collected minimizers
[M::mm_idx_gen::0.002*2.57] sorted minimizers
[M::main::0.003*2.22] loaded/built the index for 1 target sequence(s)
[M::mm_mapopt_update::0.003*2.28] mid_occ = 2
[M::mm_idx_stat] kmer size: 15; skip: 10; is_hpc: 0; #seq: 1
[M::mm_idx_stat::0.003*2.18] distinct minimizers: 74 (100.00% are singletons); average occurrences: 1.000; average spacing: 5.270
[M::worker_pipeline::0.037*2.64] mapped 1000 sequences
[M::main] Version: 2.16-r922
[M::main] CMD: minimap2 -ax map-ont Amplicon26.fasta All_genomes.fsa
[M::main] Real time: 0.037 sec; CPU: 0.098 sec; Peak RSS: 0.003 GB
(base) bas:JOVE bom86$ samtools sort All_genomes.bam > All_genomes_sorted.bam
(base) bas:JOVE bom86$ samtools stats All_genomes_sorted.bam > stats.txt
(base) bas:JOVE bom86$ ▮
```

```
(base) bas:JOVE bom86$ grep ^SN stats.txt | cut -f 2-
raw total sequences:    1000
filtered sequences:     0
sequences:      1000
is sorted:      1
1st fragments: 1000
last fragments: 0
reads mapped:   1000
reads mapped and paired:        0       # paired-end technology bit set + both mates mapped
reads unmapped: 0
reads properly paired: 0        # proper-pair bit set
reads paired:   0       # paired-end technology bit set
reads duplicated:       0       # PCR or optical duplicate bit set
reads MQ0:      0       # mapped and MQ=0
reads QC failed:        0
non-primary alignments: 0
total length:   390246  # ignores clipping
total first fragment length:    390246  # ignores clipping
total last fragment length:     0       # ignores clipping
bases mapped:   390246  # ignores clipping
bases mapped (cigar):   390246  # more accurate
bases trimmed:  0
bases duplicated:       0
mismatches:     10241   # from NM fields
error rate:     2.624242e-02    # mismatches / bases mapped (cigar)
average length: 390
average first fragment length:  390
average last fragment length:   0
maximum length: 391
maximum first fragment length:  0
maximum last fragment length:   0
average quality:        255.0
insert size average:    0.0
insert size standard deviation: 0.0
inward oriented pairs:  0
outward oriented pairs: 0
pairs with other orientation:   0
pairs on different chromosomes: 0
percentage of properly paired reads (%):        0.0
(base) bas:JOVE bom86$ gre
```

```
(base) bas:JOVE bom86$ grep ^IC stats.txt | cut -f 2-
188     0       8       0       0
191     0       238     0       0
(base) bas:JOVE bom86$
```

Suppl File 1-Script for the random selection

Click here to access/download
**Supplemental Coding Files**
Random_selection.py