

<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

# 1 MinION™ nanopore sequencing of environmental metagenomes: a synthetic approach

2  
3  
4 3 Bonnie L. Brown, Virginia Commonwealth University, Department of Biology, 1000 W Cary Street,

5  
6 4 Richmond, VA 23284, USA, [blbrown@vcu.edu](mailto:blbrown@vcu.edu)

7  
8 5 Mick Watson, The Roslin Institute, University of Edinburgh, Division of Genetics and Genomics,

9  
10 6 Easter Bush, Midlothian, EH25 9RG, UK, [mick.watson@roslin.ed.ac.uk](mailto:mick.watson@roslin.ed.ac.uk)

11  
12 7 Samuel S. Minot, One Codex, 165 11<sup>th</sup> St, San Francisco, CA 94103, USA, [sam@onecodex.com](mailto:sam@onecodex.com)

13  
14 8 Maria C. Rivera, Virginia Commonwealth University, Department of Biology, 1000 W Cary Street,

15  
16 9 Richmond, Virginia 23284, USA, [mcrivera@vcu.edu](mailto:mcrivera@vcu.edu)

17  
18 10 Rima B. Franklin, Virginia Commonwealth University, Department of Biology, 1000 W Cary Street,

19  
20 11 Richmond, Virginia 23284, USA, [rbfranklin@vcu.edu](mailto:rbfranklin@vcu.edu)

21  
22 12  
23  
24  
25  
26 13 **Corresponding author:** Bonnie L. Brown, [blbrown@vcu.edu](mailto:blbrown@vcu.edu)

## 27 28 29 30 31 15 **Abstract**

32  
33 16 **Background:** Environmental metagenomic analysis is typically accomplished by assigning

34  
35 17 taxonomy and/or function from whole genome sequencing (WGS) or 16S amplicon sequences.

36  
37 18 Both of these approaches are limited, however, by read length, among other technical and

38  
39 19 biological factors. A nanopore-based sequencing platform, MinION™, produces reads that are

40  
41 20  $\geq 1 \times 10^4$  bp in length, potentially providing for more precise assignment, thereby alleviating some

42  
43 21 of the limitations inherent in determining metagenome composition from short reads. We tested

44  
45 22 the ability of sequence data produced by MinION (R7.3 flow cells) to correctly assign taxonomy

46  
47 23 in single bacterial species runs and in three types of low complexity synthetic communities: a

48  
49 24 mixture of DNA using equal mass from four species, a community with one relatively rare (1%)

50  
51 25 and three abundant (33% each) components, and a mixture of genomic DNA from 20 bacterial

52  
53 26 strains of staggered representation. Taxonomic composition of the low-complexity communities

54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

27 was assessed by analyzing the MinION sequence data with three different bioinformatic  
28 approaches : Kraken, MG-RAST, and One Codex.

29  
30 **Results:** Long read sequences generated from libraries prepared from single strains using the  
31 version 5 kit and chemistry, run on the original MinION device, yielded as few as 224 to as many  
32 as 3,497 bidirectional high-quality (2D) reads with an average overall study length of 6,000 bp.  
33 For the single-strain analyses, assignment of reads to the correct genus by different methods  
34 ranged from 53.1% to 99.5%, assignment to the correct species ranged from 23.9% to 99.5%, and  
35 the majority of mis-assigned reads were to closely related organisms. A synthetic metagenome  
36 sequenced with the same setup yielded 714 high quality 2D reads of approximately 5,500 bp that  
37 were up to 98% correctly assigned to the species level. Synthetic metagenome ~~from~~ MinION  
38 libraries generated using version 6 kit and chemistry yielded 899-3,497 2D reads with lengths  
39 averaging 5,700 bp with up to 98% assignment accuracy at the species-level. The observed  
40 community proportions for “equal” and “rare” synthetic libraries were close to the known  
41 proportions, deviating from 0.1 – 10% across all tests. For a 20-species mock community with  
42 staggered contributions, a sequencing run detected all but 3 species (each included at <0.05% of  
43 DNA in the total mixture); 91% of reads were assigned to the correct species, 93% of reads were  
44 assigned to the correct genus, and >99% of reads were assigned to the correct family.

45  
46 **Conclusions:** At the current level of output and sequence quality (just under  $4 \times 10^3$  2D reads for a  
47 synthetic metagenome), MinION sequencing followed by Kraken or One Codex analysis has the  
48 potential to provide rapid and accurate metagenomic analysis where the consortium is comprised  
49 of a limited number of taxa. Important considerations noted in this study included: high  
50 sensitivity of the MinION platform to the quality of input DNA, high variability of sequencing  
51 results across libraries and flow cells, and relatively small numbers of 2D reads per analysis limit.  
52 Together, these limited detection of very rare components of the microbial consortia, and would  
53 likely limit the utility of MinION for the sequencing of high-complexity metagenomic

1 54 communities where thousands of taxa are expected. Furthermore, the limitations of the currently  
2 55 available data analysis tools suggest there is considerable room for improvement in the analytical  
3  
4 56 approaches for the characterization of microbial communities using long reads. Nevertheless, the  
5  
6 57 fact that the accurate taxonomic assignment of high quality reads generated by MinION is  
7  
8 58 approaching 99.5% and, in most cases, the inferred community structure mirrors the known  
9  
10 59 proportions of a synthetic mixture, warrants further exploration of practical application to  
11  
12 60 environmental metagenomics as the platform continues to develop and improve. With further  
13  
14 61 improvement in sequence throughput and error rate reduction, this platform shows great promise  
15  
16 62 for precise real-time analysis of the composition and structure of more complex microbial  
17  
18 63 communities.  
19  
20  
21

22 64

#### 23 65 **Keywords**

24 66 MinION™, Oxford Nanopore Technologies, metagenome, whole-genome sequencing, long-read  
25  
26 67 sequencing  
27  
28  
29

30 68

#### 31 69 **Background**

32  
33 70 Environmental metagenomics, employing whole genome sequence analysis to identify ecologically  
34  
35 71 and epidemiologically important components of sediments, soils, waters, and surfaces, is rapidly  
36  
37 72 evolving through advances in both hardware and software [1]. Knowledge of the consortia that  
38  
39 73 inhabit these ecosystems allows for better understanding of the organisms and their ecological roles,  
40  
41 74 provides for the development of effective strategies to mitigate ecosystem damage, and facilitates  
42  
43 75 evaluation of the responses of species to environmental change. One common approach in  
44  
45 76 environmental metagenomics involves sequencing and subsequent annotation of whole genome  
46  
47 77 nucleic acid fragments (WGS) extracted directly from environmental samples to discover major  
48  
49 78 microbial members of the ecosystem; if sequenced deeply enough, rare species can be detected [2].  
50  
51 79 For well-studied members of the microbial community, such metagenomic data also can be used to  
52  
53 80 characterize the functional potential of complex communities.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 81 One technique for characterizing environmental metagenomes is to use short-read high-  
2 82 throughput sequencing followed by mapping the reads to reference genomes. Profiling the taxonomic  
3  
4 83 composition of the community also can be accomplished by the analysis of the distribution of k-mers  
5  
6 84 (e.g., using Kraken or One Codex). Although these methodologies are very powerful due to the depth  
7  
8 85 of sequencing, the capacity to resolve the taxonomy of the community to the species level is limited  
9  
10 86 by read length. One approach to overcome this limitation is to assemble short reads into contigs prior  
11  
12 87 to analysis and annotation. If assembled correctly, the longer sequence lengths of the contigs have a  
13  
14 88 greater chance of accurately identifying the members of the community; however, due to the mixed  
15  
16 89 nature of the samples, such assembly approaches are challenged by many artifacts including chimeric  
17  
18 90 contigs that inappropriately combine sequence reads from multiple species. The high information  
19  
20 91 content of very long reads such as those provided by MinION™ (Oxford Nanopore Technologies,  
21  
22 92 Inc., Oxford, UK) has the potential to overcome some of the limitations of short reads by allowing for  
23  
24 93 longer alignments that potentially can contribute to higher taxonomic specificity, functional  
25  
26 94 characterization, and resolution. Although conceived almost two decades ago [3], nanopore-based  
27  
28 95 whole-molecule sequencing has only recently become available to MinION™ Access Programme  
29  
30 96 (MAP) participants for exploration and practical application [4]. Data generated by early access  
31  
32 97 MinION™ flow cells have been assessed for whole genome sequencing [5, 6, 7, 8, 9], gene  
33  
34 98 expression and transcriptome studies [10, 11, 12], clinical applications such as inferring antibiotic  
35  
36 99 resistance of bacterial strains and the detection of influenza and Ebola virus [13, 14, 15], bacterial  
37  
38 100 and viral serotyping [16], and clinical metagenomes of viral pathogens [17]. Efforts to use this  
39  
40 101 technology to study diverse environmental communities have been limited [18] and there has not  
41  
42 102 been, to our knowledge, any cross-validation of the results or any systematic assessment to determine  
43  
44 103 the best data analysis strategies for nanopore-based environmental metagenomics. To investigate the  
45  
46 104 potential of this platform for broader applications, we performed a set of experiments to quantify the  
47  
48 105 ability of MinION™ long-read sequence data to accurately characterize the taxonomic composition  
49  
50 106 and structure of metagenomes by assessing its performance in the characterization of low complexity  
51  
52 107 synthetic metagenomes.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

108

1  
2 **109 Data description**  
3

4 110 The raw MinION data [19] collected during sequencing by MinKNOW software (versions 0.49.2.9  
5  
6 111 through 0.51.3.40 b201605171140) were immediately uploaded as FAST5 packets to Metrichor  
7  
8  
9 112 Agent (r7.3 2D basecalling, ver rx-2.22-44717-dg-1.6.1-ch-1.6.3; Mk1 2D base-calling, ver WIMP  
10  
11 113 Bacteria k24 for SQK-MAP006), after which base-called data [19] were returned to the host  
12  
13 114 computer, also in the form of FAST5 files. The programs poRe [20], Poretools [21], and NanoOK  
14  
15 115 [22] were used to extract and characterize the numbers of reads and channels, after which only the 2D  
16  
17  
18 116 reads were stored in FASTQ and FASTA files for downstream analyses. The base-called data sets  
19  
20 117 were scrutinized by methods commonly employed in metagenome analysis of short reads including  
21  
22 118 MG-RAST [23], which assigns taxonomy based on predicted proteins and rRNA genes. The data sets  
23  
24 119 also were analyzed by newer tools aimed at long-read data including: (1) WIMP [24], which assigns  
25  
26 120 taxonomy by comparing read sequences against a database of bacteria, (2) Kraken [25], which uses  
27  
28  
29 121 exact alignments of *k-mers* and indexes more than 5000 genomes and plasmids, (3) One Codex [26],  
30  
31 122 which uses exact *k-mer* alignment to classify sequences against a reference database of ~40,000  
32  
33 123 complete microbial genomes (including bacteria, viruses, fungi, protists, and archaea), and (4) by  
34  
35 124 principal components analysis (PCA) based on the frequency of 5-mers in each read followed by  
36  
37  
38 125 annotation of reads with the top BlastN [27] hit (carried out in R [28]). Specific parameters are  
39  
40 126 described in Methods.

41  
42 127  
43  
44 **128 Results**  
45

46 129 MinION™ WGS libraries were generated from 1 µg of fresh DNA isolates (see Methods) of separate  
47  
48  
49 130 cultures of two Proteobacteria, *Escherichia coli* and *Pseudomonas fluorescens*, and two  
50  
51 131 Cyanobacteria, *Microcystis aeruginosa*, and *Synechococcus elongatus*, and from two different DNA  
52  
53 132 mixtures of these four species. One mixture combined an equal mass of genomic DNA (gDNA) from  
54  
55  
56 133 each of the four species. The other mixture was created by combining 33% mass of gDNA from each  
57  
58 134 of three species and only 1% of gDNA mass from the other species. The preparation of these libraries  
59  
60  
61  
62  
63  
64  
65

135 yielded sufficient Pre-sequencing Mix for multiple loads of each flow cell. An additional library was  
136 derived from a commercially prepared 20-species mock community. Because only 100 ng of material  
137 was provided by the supplier, genome pre-amplification using  $\Phi$ 29 polymerase was required to  
138 generate sufficient mass of DNA to create the sequencing library (see Methods).

139 To assess the purity of the cultures used in this study, we used the Sanger method to sequence  
140 full-length (~1500 bp) 16S amplicons from each (Table 1). Inspection of those data revealed varying  
141 degrees of genomic uniqueness at the species level. For the strain of *M. aeruginosa* used in this  
142 study, the top 16S hit had a low sequence identity to any reference sequence in the database (90%).  
143 In contrast, the input strain of *S. elongatus* was 99% identical to two different species of  
144 *Synechococcus* (*S. elongatus* and *S. UTEX 2973*). In addition, whole-genome alignment indicated  
145 that the input strain of *P. fluorescens* was highly similar to multiple species of *Pseudomonas*.  
146 However, all of the input organisms were distinct at the genus level, thus that taxonomic level was  
147 used for downstream analysis of the single-species and ‘Equal’ and ‘Rare’ synthetic samples.

148 MinION sequencing of the single-species libraries generated up to  $31 \times 10^3$  reads ( $0.2$ – $1.1 \times 10^3$  2D  
149 reads that passed the quality filter) ranging from as short as 5 bp to as long as  $267 \times 10^3$  bp (data  
150 include both 2D pass and fail reads), and the resulting average length of single-species read subjected  
151 to downstream analysis was  $6 \times 10^3$  bp. Using MG-RAST, Kraken, and One Codex, up to 99.5% of  
152 the high quality 2D reads obtained from the sequencing of the single-species libraries of *E. coli*, *P.*  
153 *fluorescens*, *S. elongatus*, and *M. aeruginosa* were taxonomically assigned to the corresponding input  
154 taxa. (Table 3). The least accurate assignments were for *M. aeruginosa*, where at best 58% of 2D  
155 reads were correctly assigned to the level of species, although more than half of the mis-assigned  
156 reads were to closely related cyanobacteria genera and other prokaryotes known to break down  
157 microcystin [29] (data not shown). All three methods of analysis assigned sequence reads of the *P.*  
158 *fluorescens* single-species library to *Stenotrophomonas*. Over all of these analyses, MG-RAST  
159 generally showed the lowest rate of correct taxonomic assignment and, although One Codex and  
160 Kraken provided similar results, Kraken showed a lower rate of correct assignment for *M. aeruginosa*  
161 (85%) compared to One Codex (95%).

162 In the second round of validation, using three synthetic communities containing mixtures of the  
163 previously described species,  $6\text{--}12\times 10^3$  reads ( $0.7\text{--}1.3\times 10^3$  2D reads) were generated per run, ranging  
164 in length from  $0.6\text{--}56.8\times 10^3$  bp (Table 2). For the two communities comprised of equal DNA  
165 contribution from four bacteria (25% each species), WGS proportions accurately aligned with the  
166 known proportions 87–99% of the time when analyzed using Kraken or One Codex and 65–85%  
167 using MG-RAST (Table 3). Specifically, taxonomic assignment of reads obtained from the  
168 sequencing of the equal mixture of four species (25% of each) using version 5 chemistry and run on  
169 an original MinION device identified the following taxa: 27% *E. coli*, 16% *M. aeruginosa*, 30% *P.*  
170 *fluorescens*, 21% *S. elongatus*, 3% Enterobacteriaceae, and 3% misclassified. In a subsequent test  
171 (version 6 chemistry), classification results for the equal mixture were: 26% *E. coli*, 18% *M.*  
172 *aeruginosa*, 30% *P. fluorescens*, 22% *S. elongatus*, and 3% Enterobacteriaceae, and 1% misclassified  
173 (Figure 1). For the community with three common (33% of each) and one rare (1%) representative,  
174 classifications were: 33% *E. coli*, 34% *P. fluorescens*, 29% *S. elongatus*, 1% *M. aeruginosa*, 2%  
175 misclassified (a third of those latter category of reads were assigned to *Shigella*). For both the  
176 “Equal” and “Rare” community data sets, the 5-mer frequency profiles were computed and visualized  
177 using the top BlastN hit for each full read, revealing that 5-mer profiles for these long-read sequences  
178 were shared within species. This was reflected in the 5-mer frequency analysis which revealed  
179 distinct clusters in the PCAs (Figure 2).

180 In the final round of testing, the mock microbial community with 20 species included in  
181 “staggered” proportions (i.e., 1,000 to 1,000,000 16S rRNA operon copies per organism per  $\mu\text{L}$  of  
182 material supplied by BEI Resources, Catalog # HM-783D) yielded  $14.7\times 10^3$  reads ( $3.5\times 10^3$  2D reads)  
183 ranging in length from  $0.5\text{--}20.9\times 10^3$  bp, sufficient to detect all of the high and moderate abundance  
184 species, but the sequencing run failed to detect 3 of 5 species that were included at very low mass  
185 ( $0.6\text{--}1.0$  pg/ $\mu\text{L}$  of material supplied; Table 4). For that run, misclassifications accounted for only  
186 0.2% of read assignments but greatly overrepresented in the results for this run were reads assigned to  
187 *E. coli* (included as 20% of DNA but observed as 46–52% of read assignments), whereas greatly  
188 underrepresented in the results were reads assigned to *R. sphaeroides*, which was putatively included



189 as 41% of DNA mass but accounted for only 1% of read assignments (Figure 3). Although 75% of  
 190 the read assignments made by WIMP were to genera known to comprise the mock community, 93%  
 191 of the read assignments made by One Codex matched the correct genera.

192  
193

**Table 1** Identity of single-species used in this study as determined by Sanger sequencing of  
 16S rDNA amplicons from different DNA preparations of each species.

Culture <sup>a</sup>	Final sequence Length (bp)	Sequence matches in BlastN	
		%	Organism
<i>Escherichia coli</i>	1440–1696 <sup>a</sup>	98	<i>E. coli</i> numerous strains
<i>Microcystis aeruginosa</i>	1418	90	<i>M. aeruginosa</i> NIES-843 and NIEHS-2549, and <i>M. panniformis</i> FACHB-1757
<i>Pseudomonas fluorescens</i>	1478–1570	96	<i>P. fluorescens</i> A506 and LBUM223
<i>Synechococcus elongatus</i>	1431–1719	99	<i>S. elongatus</i> PCC 7942, PCC 6301, UTEX 2973

<sup>a</sup> Multiple DNA preparations from bacterial cultures were used during the progress of the study, and each was tested,  
 yielding for each strain slightly different final 16S sequence lengths, but the same BLAST matches.

198  
199

**Table 2** Details of MinION™ whole genome sequencing output for single-species and  
 synthetic mixtures. Sequencing experiments used the MinION device and new R7.3  
 flow cells. Libraries were prepared with kit SQK–MAP005 as indicated by (5) and  
 SQK–MAP006 chemistry, indicated by (6). Columns relating to “2D” indicate bi-  
 directional reads with quality above Q9.

Experiment (chemistry)	Pores with reads	Run time (hr) <sup>a</sup>	Total bp (Mbp)	Total reads	Number of 2D pass reads	Mean 2D read length (bp)	MG-RAST accession	ENA accession
Single species								
<i>E. coli</i> <sup>(5)</sup>	430	42	83.6	26590	1112	5274	4629367.3	ERR1713483
<i>P. fluorescens</i> <sup>(5)</sup>	453	48	119.4	25228	777	7784	4629445.3	ERR1713487

<i>M. aeruginosa</i> <sup>(5)</sup>	377	18	40.8	22760	569	5676	4629369.3	ERR1713486
<i>S. elongatus</i> <sup>(5)</sup>	367	23	18.3	6163	224	5101	4629381.3	ERR1713489
Mixtures								
Equal <sup>(5)</sup>	129	24	26.5	10592	714	5527	4614572.3	ERR1713484
Equal <sup>(6)</sup>	437	44	77.1	12174	1358	5202	4685746.3	ERR1713485
Rare <sup>(6)</sup>	449	18	39.0	6728	899	6194	4685745.3	ERR1713488
Staggered <sup>(6)</sup>	300	33	39.0	14711	3497	2612	4705090.3	ERR1713490

<sup>a</sup> Runs were set to either 24 or 48 hours and were allowed to continue until either sufficient sequence data were collected or until the 2D pass rate was greatly reduced.

**Table 3** Taxonomic assignment accuracy of metagenomic reads across three analysis methods.

Accuracy of assignment to known genus (%)			
Experiment	MG-RAST	Kraken	One Codex
Single species			
<i>E. coli</i> <sup>(5)</sup>	74.4 <sup>a</sup>	99.5	98.7
<i>P. fluorescens</i> <sup>(5)</sup>	84.9 <sup>b</sup>	84.6 <sup>b</sup>	84.2 <sup>b</sup>
<i>M. aeruginosa</i> <sup>(5)</sup>	53.1	85.8	95.1
<i>S. elongatus</i> <sup>(5)</sup>	87.9	98.1	97.6
Mixtures			
Equal <sup>(5)</sup>	65.0 <sup>b</sup>	97.6	87.4 <sup>c</sup>
Equal <sup>(6)</sup>	85.9	98.0	98.7
Rare <sup>(6)</sup>	92.9	99.1	98.7

<sup>a</sup> 15% of reads assigned to *Shigella*

<sup>b</sup> 7-15% of reads assigned to *Stenotrophomonas*

<sup>c</sup> 7% of reads assigned to *Stenotrophomonas*

Accuracy was calculated as the proportion of reads assigned to the known input organism at the genus level out of the total number reads given any assignment at that rank.

**Table 4** Known composition of 20-species mock staggered community compared with analysis results for WIMP and One Codex. “nd”: not detected.

Organism	Operon count/μL <sup>a</sup>	Quantity pg/μL <sup>b</sup>	% DNA in template <sup>c</sup>	WIMP % species	WIMP % genus	One Codex % species	One Codex % genus
----------	------------------------------	-----------------------------	--------------------------------	----------------	--------------	---------------------	-------------------

<i>Acinetobacter baumannii</i>	10,000	8.2	0.24	0.14	0.14	0.29	0.29
<i>Actinomyces odontolyticus</i>	1,000	1	0.03	nd	nd	nd	nd
<i>Bacillus cereus</i>	100,000	45	1.33	0.53	0.53	0.66	0.75
<i>Bacteroides vulgatus</i>	1,000	0.8	0.02	0.1	0.1	0.07	0.12
<i>Clostridium beijerinckii</i>	100,000	44	1.30	0.19	0.19	0.29	0.35
<i>Deinococcus radiodurans</i>	1,000	1	0.03	0.05	0.05	0.07	0.06
<i>Enterococcus faecalis</i>	1,000	0.7	0.02	nd	nd	nd	nd
<i>Escherichia coli</i>	1,000,000	680	20.04	45.61	45.66	52.15	52.52
<i>Helicobacter pylori</i>	10,000	8.6	0.25	1.68	1.68	3.43	2.72
<i>Lactobacillus gasseri</i>	10,000	3.2	0.09	0.14	0.14	0.22	0.23
<i>Listeria monocytogenes</i>	10,000	5	0.15	0.38	0.38	0.58	0.52
<i>Neisseria meningitidis</i>	10,000	5.8	0.17	0.24	0.24	0.44	0.41
<i>Propionibacterium acnes</i>	10,000	8.8	0.26	0.48	0.48	0.07	0.64
<i>Pseudomonas aeruginosa</i>	100,000	160	4.71	1.25	1.25	3.07	3.18
<i>Rhodobacter sphaeroides</i>	1,000,000	1,400	41.25	1.01	1.01	1.46	1.27
<i>Staphylococcus aureus</i>	100,000	59	1.74	0.38	3.88	1.31	12.74
<i>Staphylococcus epidermidis</i>	1,000,000	510	15.03	7.67	7.72	6.65	
<i>Streptococcus agalactiae</i>	100,000	32	0.94	0.96	1.01	0.95	16.97
<i>Streptococcus mutans</i>	1,000,000	420	12.38	10.17	10.17	19.50	
<i>Streptococcus pneumoniae</i>	1,000	0.6	0.02	nd	nd	nd	
Other		0	0	29.02 <sup>d</sup>	25.37 <sup>e</sup>	8.77 <sup>f</sup>	7.24 <sup>g</sup>
Correct assignments				70.98	74.63	91.23	92.76

<sup>a</sup>Theoretical copy number provided by BEI Resources certificate of analysis

<sup>b</sup>gDNA content provided by BEI Resources certificate of analysis

<sup>c</sup>Proportion of individual species within the mock community.

<sup>d</sup>Of these, 12.7% were correctly assigned to genus, 86.4% were Enterobacteriaceae, and only 0.7% were misclassifications.

<sup>e</sup>Of these, 86.4% were Enterobacteriaceae and only 0.7% were misclassified.

<sup>f</sup>Of these, 56.8% were *Shigella*.

<sup>g</sup>Of these, 63.3% were species of *Escherichia* and *Shigella*.

## Discussion

Sequencing of whole genome libraries can enhance environmental metagenomic analysis by providing more precise identification of the composition and structure of the community than is

228 possible by amplicon sequencing of marker genes (e.g., 16S) [2, 30]. Typical environmental samples  
1  
2 229 contain tens of thousands to millions of organisms, yet the resulting metagenomes almost certainly  
3  
4 230 underrepresent this diversity and, often due to short-read strategy, the resulting data sets can be  
5  
6 231 confidently assigned only to higher taxonomic levels [31, 32]. One strategy to improve the accuracy  
7  
8 232 of the taxonomic assignment is to carefully assemble metagenomic data, which despite the resulting  
9  
10 233 chimerism has been shown to greatly enhance species call correctness [33]. However, even with  
11  
12 234 enhanced sequencing and bioinformatic strategies, many public database accessions contain  
13  
14 235 sequences that are not innate to the species that was analyzed; these include symbionts, parasites,  
15  
16 236 pathogens, and sequencing linkers/primers/adapters (unbeknownst to those who have accessed the  
17  
18 237 data) that can lead to false discovery rates [34]. Contaminated and mis-annotated reference sequences  
19  
20 238 can affect environmental metagenome analyses that are derived from short reads to a greater extent  
21  
22 239 than would be expected from analyses based on long reads. Long reads can circumvent these issues  
23  
24 240 [31, 35, 36], so long as much of the genome for each component organism is represented in the  
25  
26 241 sequencing library and there are few errors in the sequences and the reference database. The results  
27  
28 242 reported here allow us to consider the potential utility of MinION long read sequencing and  
29  
30 243 subsequent bioinformatic analysis for shotgun environmental metagenomics.

31 244 The primary challenge of microbial metagenomic sequence analysis using long reads is the  
32  
33 245 comparison of input sequences against a large reference database of whole genomes from bacteria,  
34  
35 246 viruses, fungi, etc. Although a number of algorithms have been developed for alignment of long,  
36  
37 247 error-prone reads [37, 38], those sensitive algorithms are not optimized for the challenge of  
38  
39 248 comparison against the large and ever-expanding universe of microbial genomes. The bioinformatic  
40  
41 249 methods used in this analysis, MG-RAST, Kraken, One Codex, and WIMP, each compare the input  
42  
43 250 reads against their own more concise reference databases, providing an assignment for the most likely  
44  
45 251 origin of each individual sequence.

46 252 We found that for low complexity synthetic communities, long reads generated by MinION  
47  
48 253 provided sufficiently precise sequence data to assign organisms represented at or above 1%. In fact,  
49  
50 254 two out of five species included at <0.05% in a mock community (and 9 out of 9 species included at  
51  
52 255 0.05-1.00%) were detected. Furthermore, for un-amplified whole genome preparations, read  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

256 assignments were observed to be within about 10% of their proportional occurrence in the  
1  
2 257 metagenome. Ultimately, we saw that although the reads were longer, because the sequence coverage  
3  
4 258 was not as deep, the improvement in specificity of assignment was offset by a reduction in the  
5  
6 259 sensitivity, and some of the genomes present at low concentration were not detected.

8  
9 260 By comparing the output of multiple analysis methods, we were able to gain insight into the  
10  
11 261 performance of various bioinformatic approaches for analyzing error-prone MinION reads. Overall,  
12  
13 262 MG-RAST provided the lowest level of accuracy and detected multiple organisms that were not a part  
14  
15 263 of the known input set. This is not surprising given that MG-RAST is optimized for analyzing short-  
16  
17 264 read, low-error data. Kraken and One Codex performed similarly for the single-species samples  
18  
19 265 except in the case of *M. aeruginosa*, in which case One Codex correctly identified this taxon at a  
20  
21 266 higher rate than Kraken (95% vs. 85%). An unexpected finding of this study was the detection by all  
22  
23 267 three methods of *Stenotrophomonas* in the *P. fluorescens* single-species sample. Interestingly,  
24  
25 268 *Stenotrophomonas* was classified as *Pseudomonas* when it was first discovered, based on similar  
26  
27 269 metabolic capabilities, and was later moved to its own genus based on molecular data [39]. Our 16S  
28  
29 270 sequences derived from laboratory cultures used in this study did not identify *Stenotrophomonas*,  
30  
31 271 suggesting that its identification in the mixed metagenomes is not a result due to a contaminant but  
32  
33 272 rather, an artifact caused by assigning taxonomy to reads with multiple sequencing errors. Also  
34  
35 273 contributing to its identification is the fact that both *Pseudomonas* and *Stenotrophomonas* share  
36  
37 274 functional phenotypic characteristics, indicating they may share homologous genes coding for those  
38  
39 275 characteristics. The sharing of homologous genes, similar GC contents (both species genomes have  
40  
41 276 66% GC), and the higher error rate are the most likely factors responsible for the assignment of  
42  
43 277 *Pseudomonas* sequence reads to *Stenotrophomonas*.

44  
45 278 The fact that the estimated proportions of community members in synthetic mixtures were not  
46  
47 279 precise despite careful DNA quantitation could indicate differences across library preparation (all  
48  
49 280 libraries were prepared by BLB), reagent kits, flow cells, MinKNOW control scripts, the quality of  
50  
51 281 DNAs used to create the synthetic metagenomes, and the methods used for quantification (Qubit for  
52  
53 282 the home-grown mixtures and UV spectrophotometry for the 20-species mixture). Because DNA  
54  
55 283 quality is of paramount importance for MinION sequencing, PreCR (used in the version 5 protocol) or  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

284 FFPE Repair Mix (used in the version 6 protocol) was included in the preparation of all libraries. The  
1  
2 285 potential for profound effects related to library preparation recently was examined by Jones and  
3  
4 286 collaborators [30], leading to the recommendation that studies of complex metagenomes should be  
5  
6 287 based on PCR-free approaches. The current data indicate that the MinION lends itself well to a PCR-  
7  
8 288 free approach but its utility for the analysis of complex metagenomes is presently limited by the small  
9  
10 289 number of reads that pass the quality filtering process. The current study also provides data for  
11  
12 290 considering alternatives to PCR for amplification, in this case GenomiPhi™, which was used to  
13  
14 291 generate sufficient DNA for one library in the current study (“Staggered”). This method is optimized  
15  
16 292 for linear DNA and was intended to generate unbiased copies of the 20-species genomes.  
17  
18 293 Nevertheless, the Φ29 pre-amplification step is one possible reason for the overrepresentation of *E.*  
19  
20 294 *coli* and underrepresentation of *R. sphaeroides* in the sequencing of the 20 species mock community.  
21  
22 295 Also, a consequence of Φ29 pre-amplification combined with putative differences in DNA quality,  
23  
24 296 chimeric amplicons (known to occur with Φ29 amplification of microbial communities [40]) could  
25  
26 297 have been formed predominantly from higher quality *E. coli* DNA re-priming itself [41] leading to  
27  
28 298 overrepresentation of the *E. coli* component. Notably, a novel low input DNA approach recently  
29  
30 299 reported [42] could enhance MinION analyses of samples with low DNA yields. Although the pre-  
31  
32 300 amplification step is the most likely culprit, an additional effect that could contribute to incongruence  
33  
34 301 of known and estimated proportions in the 20 species mock community is that organisms for which  
35  
36 302 there are many accessions in the public databases provide for more precise classification (e.g., NCBI  
37  
38 303 has more than  $6 \times 10^5$  *E. coli* complete genome accessions) and that *vice versa*, organisms with  
39  
40 304 relatively few accessions (e.g., NCBI has only 116 *R. sphaeroides* complete genome accessions)  
41  
42 305 result in less precise classification.  
43  
44  
45  
46  
47  
48

49 306 Despite the rather small number of 2D reads that were observed to pass the quality filter across all  
50  
51 307 MinION runs, there was a strong biological signal in the data (Figure 2). Thus, as investigators have  
52  
53 308 found MinION useful for single genome introspection [6, 9, 15], 16S and other amplicon resolution  
54  
55 309 [16, 43], cDNA sequencing [11], and assembly [5, 44, 45], our findings imply that this platform has  
56  
57 310 immediate utility for analysis of very simple mixtures (e.g., serum testing for pathogens). Over the  
58  
59  
60  
61  
62  
63  
64  
65

18-month period of MinION use for this set of experiments, 2D pass rates increased from 2% to 24%. Because the rate of improvement is concurrent with Moore's Law [46], we speculate that future improvements will make the MinION platform very useful in the analysis of complex metagenomic samples in the near future. The cloud-based WIMP base-calling and taxon prediction program associated with the device provides a method of real-time analysis of metagenomic data. However, because we had no control over the comparative database, the cloud implementation of WIMP was far less flexible for environmental metagenomic analysis than Kraken or One Codex, and we note that use of an incomplete database can lead to false positives and negatives. By the time of submission of this study, the R7.3 flow cells and sequencing chemistry were no longer available. Subsequent versions of the platform have shown dramatically lower error and higher throughput. This study nevertheless provides a baseline for considering nanopore metagenomics and provides an impetus for further development of MinION output and data analysis, specifically with regard to evaluation of the informative value of 1D reads, scrutiny of reference data, alternative alignment algorithms, and more sophisticated k-mer analyses. As the quality rate for this platform improves, the potential will increase for MinION to accurately resolve the diversity and composition of many of the taxa in an environmental metagenome.

327

## 328 **Methods**

To set a baseline of expectations for MinION metagenomic analysis, we performed single-species sequencing runs with four organisms. Cell cultures at log phase were harvested by spinning 15 mL culture tubes at  $3,000 \times g$  for 30 min, and DNA was isolated using the PowerSoil DNA kit (MoBio, Carlsbad, CA, USA) according to the manufacturer's instructions. Nucleic acid quality and quantity were checked via Nanodrop 2000 and Qubit, whereafter 1  $\mu\text{g}$  of DNA was used to prepare sequencing libraries. For the first two mixtures, equal portions of DNAs from all four organisms (250 ng each) were used ("equal") and, for the third mixture ("rare"), equivalent amounts of three of the species were used (330 ng each) and *M. aeruginosa* was included as only 1% of the mixture (10 ng). An additional preparation of a mock community containing DNA of 20 bacterial species in staggered

338 amounts was obtained from a commercial source (Catalog # HM-783D, BEI Resources, ATCC,  
1  
2 339 Manassas, VA, USA). This mock community preparation was chosen because it previously has been  
3  
4 340 used to test the ability of the R7.3 version MinION to study microbial diversity via 16S amplicon  
5  
6 341 approach [43]. However, because sequencing libraries for this study required 1 µg of DNA, to  
7  
8  
9 342 generate sufficient starting material 1 µL of the mock community sample (5.5 ng of template, the  
10  
11 343 amount recommended by the supplier for a typical reaction) was pre-amplified using Φ29 enzyme  
12  
13 344 from the GenomiPhi V3 kit (25-6601-24, GE Healthcare Bio-Sciences, Pittsburgh, PA, USA)  
14  
15 345 according to the manufacturer's recommendations. This version of Φ29 enzyme was chosen for  
16  
17  
18 346 isothermal pre-amplification due to the high-fidelity proof-reading aspects of its replication process  
19  
20  
21 347 [47].

22  
23 348 The composition of each microbial mixture was calculated on the basis of the relative DNA mass  
24  
25 349 contributed from each organism. Due to the random nature of shotgun sequencing, this library  
26  
27 350 construction strategy is expected to result in a relative proportion of reads sequenced from each  
28  
29  
30 351 organism that corresponds to the relative input mass. In other words, the relative genome size of each  
31  
32 352 organism should not have impacted the relative proportion of reads recovered from each organism.

33  
34 353 Sequencing libraries were prepared for R7.3 flow cells run on an original MinION device using  
35  
36 354 the Genomic DNA Sequencing Kit SQK-MAP005 (version 5 chemistry) according to the base  
37  
38 355 protocol from Oxford Nanopore with slight modifications [48] and for flow cells run using the  
39  
40  
41 356 Nanopore Sequencing Kit SQK-MAP006 (version 6 chemistry) according to the manufacturer's  
42  
43 357 recommendations. The steps for library SQK-MAP005 preparation included in this order: shearing 1  
44  
45 358 µg in a Covaris g-TUBE (Covaris, Inc., Woburn, MA, USA) at 2,000 × g for 2 min, treatment with  
46  
47 359 PreCR (New England Biolabs, Beverly, MA, USA), cleanup with 1X AMPure beads (Agencourt,  
48  
49  
50 360 Beckman Coulter, Brea CA, USA), end-repair with NEBNext End Repair Module (New England  
51  
52 361 Biolabs), cleanup with 0.5X AMPure beads, dA-tailing with NEBNext dA-Tailing Module (New  
53  
54 362 England Biolabs), ligation to a cocktail of both the leader and hairpin sequencing adapters (Oxford  
55  
56 363 Nanopore Technologies) using Blunt TA Ligase (New England Biolabs), cleanup using his-tag  
57  
58  
59 364 Dynabeads (Life Technologies, City, State, USA), and recovery of the pre-sequencing mix in 25 µL  
60  
61  
62  
63  
64  
65



1 365 of Elution Buffer (Oxford Nanopore Technologies). After priming the flow cell with EP solution  
2 366 according to the manufacturer's recommendations, an initial 6  $\mu$ L aliquot of the pre-sequencing mix  
3  
4 367 (at 10-20 ng/ $\mu$ L) was combined with 141  $\mu$ L EP Solution and 3  $\mu$ L Fuel Mix and applied to the flow  
5  
6  
7 368 cell. Thereafter, at 6-8 hr intervals, additional pre-sequencing mix aliquots (held on ice) combined  
8  
9 369 with EP Solution and Fuel Mix were added to the flow cell at times roughly coinciding with re-  
10  
11 370 programmed pore "remux," which is a process that adjusts the bias voltage and mux channels to  
12  
13 371 maximize yield performance. Modified scripts (J. Tyson, pers. comm.) caused the MinION device to  
14  
15  
16 372 perform four remux steps at 8 h intervals to maintain regular increases in data (Figure 4).

17  
18 373 Steps for library SQK-MAP006 preparation included in this order: shearing in a Covaris g-TUBE  
19  
20 374 (Covaris, Inc., Woburn, MA, USA) at 2,000  $\times$  g for 2 min, treatment with PreCR (New England  
21  
22 375 Biolabs, Beverly, MA, USA), cleanup with 1X AMPure beads (Agencourt, Beckman Coulter, Brea  
23  
24  
25 376 CA, USA), combined end-repair and dA-tailing with NEBNext UltraII End Repair/dA-Tailing  
26  
27 377 Module (New England Biolabs), cleanup with 1X AMPure beads, ligation to a cocktail of both the  
28  
29 378 leader and hairpin sequencing adapters (Oxford Nanopore Technologies) using Blunt TA Ligase  
30  
31 379 (New England Biolabs), addition of a tether to the hairpin segment, cleanup using MyOne  
32  
33  
34 380 Streptavidin C1 Beads (Life Technologies, Carlsbad, CA, USA), and recovery of the pre-sequencing  
35  
36 381 mix in 25  $\mu$ L of Elution Buffer (Oxford Nanopore Technologies). After priming the flow cell with  
37  
38 382 running buffer and fuel according to the manufacturer's recommendations, an initial 6  $\mu$ L aliquot of  
39  
40  
41 383 the pre-sequencing mix (at 10-20 ng/ $\mu$ L) was combined with 75  $\mu$ L Running Buffer, 65  $\mu$ L water, and  
42  
43 384 4  $\mu$ L Fuel Mix and applied to the flow cell. Thereafter, at 8 hr intervals, additional pre-sequencing  
44  
45 385 mix aliquots (held on ice) were combined with Running Buffer and Fuel Mix and added to the flow  
46  
47  
48 386 cell at times roughly coinciding with re-programmed pore remux (modified scripts from J. Tyson,  
49  
50 387 pers. comm.) Modified remux scripts were not used for the final MinION run (Staggered community  
51  
52 388 analysis) because that run was controlled by a new version of MinKNOW.

53  
54 389 Whole genome sequence data (2D FASTQ) from the MinION R7.3 flow cells were accessed on  
55  
56  
57 390 the MG-RAST server [23] and annotated based on their predicted proteins and rRNA genes using the  
58  
59 391 BLAT annotation algorithm [49] against the M5NR protein Db, screened to remove any sequences  
60  
61  
62  
63  
64  
65

392 matching *H. sapiens* (none found), and without dereplication or dynamic trimming. Although  
1  
2 393 optimized for short read data, the MG-RAST tools were implemented because they allow query of a  
3  
4 394 suite of comprehensive nonredundant genetic databases and because this server provides a means to  
5  
6 395 share both raw data and computational results. Raw read counts were later accessed from MG-RAST  
7  
8 396 using the API endpoint for organism summaries. The recommended parameters “hit\_type=single”,  
9  
10 397 “source=RefSeq”, and “evaluate=15” were used to generate the appropriate read-level abundance  
11  
12 398 information. The same read sets (2D FASTA) also were analyzed by Kraken [25] using the default k-  
13  
14 399 mer size, minimizers, and other parameters, and accessing a local database created from archaea,  
15  
16 400 bacteria, fungi, virus, protozoa, human, and invertebrate genomes. The Kraken tool was implemented  
17  
18 401 because it is much faster than MG-RAST and allowed use of a smaller, more targeted reference  
19  
20 402 database. The results were translated (kraken-translate) and summarized (kraken-report) to provide  
21  
22 403 full taxonomic names for each classified sequence. Metagenomic analysis using One Codex was  
23  
24 404 performed by uploading the 2D FASTQ data to the One Codex platform at <https://app.onecodex.com>.  
25  
26 405 This cloud-based k-mer method was selected because it is reportedly more accurate than either the  
27  
28 406 MG-RAST or the Kraken tools and because like MG-RAST, it provides for community access to the  
29  
30 407 data and analytical results. Because of the high error rate of the R7.3 version MinION nucleotide  
31  
32 408 data, the unfiltered One Codex results were used for this analysis, which do not include an automated  
33  
34 409 error-filtering step. The One Codex read-level classification results were accessed by selecting the  
35  
36 410 “unfiltered” option in the web-based results display and downloading a data table for each sample to  
37  
38 411 generate appropriate read-level abundance information for tabulation.  
39  
40  
41  
42  
43

44 412 Comparative data sets were generated for each of the four single species templates using full  
45  
46 413 length ~1500 bp Sanger sequencing of a 16S amplicons [50]. Reads from the 16S analysis were  
47  
48 414 subjected to BlastN for taxonomic assignment.  
49  
50

51 415  
52  
53  
54 416 **Figure 1 Result of “What’s in my pot” analysis of a mixture with equal DNA mass from**  
55  
56 417 **four bacterial strains.**  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 418 Rendering of real-time analysis using WIMP [20] of whole genome sequences from a  
2 419 synthetic mixture prepared from equal DNA quantities of four cultured microbe species  
3  
4 420 (experiment “Equal<sup>(6)</sup>” in Tables 1 and 2) and run on the MinION™ sequencing platform.  
5  
6  
7 421 Arc angle is proportional to the number of reads assigned to the indicated species. Colors  
8  
9 422 (scale at bottom of diagram) refer to the classification score threshold (for this analysis the  
10  
11  
12 423 threshold for inclusion was 0.01).

13  
14  
15 424

16 425 **Figure 2 Principal component analysis of normalized 5-mer frequency (i.e., percentage)**  
17  
18  
19 426 **within each MinION™ read for a mixture with equal DNA mass from four bacterial**  
20  
21 427 **strains and a mixture with one rare component.**

22  
23  
24 428 A: Sequencing run with equal DNA mass from four species. B: sequencing run with three  
25  
26 429 equally represented (33% DNA mass each) and one rare (1% DNA mass) species included in  
27  
28  
29 430 the DNA pool. “none”: read had no BlastN hits. “other”: read had BlastN hits but not one of  
30  
31 431 the four species included in the mix.

32  
33  
34 432

35 433 **Figure 3 Log abundance of reads assigned from staggered mixture.**

36  
37  
38 434 DNA of 20 species mixed in various proportions (BEI Resources, ATCC, HM-783D, operon  
39  
40 435 counts  $\mu\text{L}^{-1}$  in original mixture indicated along bottom margin of bars) was pre-amplified  
41  
42  
43 436 with  $\Phi 29$  polymerase prior to library preparation and sequenced with MinION™ R7.3 flow  
44  
45 437 cells. The 2D reads that passed quality filtering were assigned to taxa using Kraken. Colored  
46  
47  
48 438 bars are species included in the mix whereas gray bars indicate species detected but not  
49  
50 439 included in the original DNA mixture.

51  
52  
53 440

54  
55 441 **Figure 4 Read production using a MinION™ device and an R7.3 flow cell.**  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

442 Illustration of reads collected from a synthetic metagenome made with equal DNA mass from  
1  
2 443 four microbias species and a library prepared using SQK-MAP006 kit. Arrows indicate  
3  
4  
5 444 approximate times when additional aliquots of library and fuel were added.  
6

7 445

#### 10 446 **Availability and requirements**

- 12 447 • Project name: Experimental Metagenome on MinION
- 13  
14  
15 448 • Project home page: <https://github.com/gigascience/paper> link will be here.
- 16  
17 449 • Operating system: Unix
- 18  
19  
20 450 • Programming language: Bash and R
- 21  
22 451 • Other requirements: Unix
- 23  
24  
25 452 • License: N/A

26  
27 453

#### 30 454 **Availability of supporting data**

32 455 The datasets supporting the results of this article are available in the GigaDB repository [19],  
33  
34 456 on the MG-RAST server 4629367.3, 4629445.3, 4629369.3, 4629381.3, 4614572.3,  
35  
36  
37 457 4685746.3, 4685745.3, 4705090.3, and at the European Nucleotide Archive as primary  
38  
39 458 accessions PRJEB8672 and PRJEB8716. One Codex results are available at  
40  
41  
42 459 [https://app.onecodex.com/projects/bb\\_minion\\_env](https://app.onecodex.com/projects/bb_minion_env).

43  
44 460

#### 47 461 **Competing interests**

49 462 BLB, MW, MCR, and RBF are enrolled in the Oxford Nanopore MinION™ Access  
50  
51  
52 463 Programme (MAP) and received free materials for this research. SSM is an employee of One  
53  
54 464 Codex.

55  
56 465

#### 59 466 **Authors' contributions**

60  
61  
62  
63  
64  
65

1 467 BLB conceived of the study, performed the DNA extraction and sequencing, directed the  
2 468 data analysis, and drafted the manuscript. MW provided bioinformatic analyses and  
3  
4 469 statistical analyses. MCR participated in study design, sequence alignment, and  
5  
6  
7 470 bioinformatic analysis. RBF participated in study design, sequencing, data analysis, and  
8  
9  
10 471 manuscript preparation. SSM performed some of the bioinformatic analyses and data  
11  
12 472 interpretation. All authors read and approved the final manuscript.  
13  
14  
15 473

## 16 474 **Acknowledgements**

17  
18  
19 475 This work was supported by the Virginia Commonwealth University Department of Biology  
20  
21  
22 476 (BLB, MCR, RBF), and by GenEco, LLC, Richmond, Virginia to BLB. Funding for MW  
23  
24 477 was from the Biotechnology and Biological Sciences Research Council including Institute  
25  
26  
27 478 Strategic Programme and National Capability grants (BBSRC; BBS/E/D/20310000,  
28  
29 479 BB/J004243/1, BB/M020037/1). The authors acknowledge M. Kensey Barker (VCU) for  
30  
31  
32 480 assistance with culturing bacteria. John Tyson (UBC) provided runtime plots and wrote the  
33  
34 481 python scripts used to control the MinION device during the run. Hugh Eaves (VCU)  
35  
36  
37 482 provided programming assistance. Sarah Highlander (Venter Inst.) provided advice on  
38  
39 483 determining DNA concentration. Michael Micorescu (ONT) provided assistance with  
40  
41  
42 484 Kraken. Arwyn Edwards and Kevin Keegan provided thorough review and suggestions for  
43  
44 485 improvement of the manuscript. The following reagent was obtained through BEI Resources,  
45  
46 486 NIAID, NIH as part of the Human Microbiome Project: Genomic DNA from Microbial  
47  
48  
49 487 Mock Community B (Staggered, Low Concentration), v5.2L, for 16S rRNA Gene  
50  
51 488 Sequencing, HM-783D.  
52  
53  
54 489

## 55 56 490 **List of abbreviations**

57  
58 491 MAP: MinION™ Access Programme  
59  
60  
61  
62  
63  
64  
65

492 2D: refers to sequences where both the template and the complement were completed

493 (bidirectional) and passed the Metrichor quality threshold (Q9)

494 PCA: principal components analysis

495 gDNA: genomic DNA isolates from putatively pure cultures of bacterial strains

496

## 497 **References**

498 1. Mendoza MLZ, T Sicheritz-Ponten, MTP Gilbert. Environmental genes and genomes:

499 understanding the differences and challenges in the approaches and software for their

500 analyses. *Briefings in Bioinformatics*, 2015, 1–14. doi: 10.1093/bib/bbv001

501 2. Thomas T, Gilbert J, Meyer F. Metagenomics – a guide from sampling to data analysis. *Microb*

502 *Inform Experim.* 2012;2:3.

503 3. Kasianowicz JJ, Brandin E, Branton D, Deamer DW. Characterization of individual polynucleotide

504 molecules using a membrane channel. *Proc Natl Acad Sci.* 1996;93:13770–73.

505 4. Eisenstein M. Oxford Nanopore announcement sets sequencing sector abuzz. *Nat Biotechnol.*

506 2012;30:295–6.

507 5. Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using only

508 nanopore sequencing data. *Nat Methods.* 2015;12:733–5.

509 6. Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz M, McCombie WR. Oxford

510 Nanopore sequencing and de novo assembly of a eukaryotic genome. *Genome Res.*

511 2015a;25:1-7.

512 7. Risse J, Thomson M, Blakely G, Koutsovoulos G, Blaxter M, Watson M. A single chromosome

513 assembly of *Bacteroides fragilis* strain BE1 from Illumina and MinION nanopore sequencing

514 data. *GigaScience* 2015;4:60.

- 515 8. Quick J, Ashton P, Calus S, Chatt C, Gossain S, et al. Rapid draft sequencing and real-time  
1 nanopore sequencing in a hospital outbreak of *Salmonella*. *Genome Biol.* 2015;16:114.  
2 516  
3  
4  
5 517 9. Madoui M-A, Engelen S, Cruaud C, Belser C, Bertrand L, et al. Genome assembly using nanopore-  
6  
7 518 guided long and error-free DNA reads. *BMC Genomics.* 2015;16:327.  
8  
9  
10 519 10. Mongan AE, Yusuf I, Wahid I, Hatta M. The evaluation on molecular techniques of reverse  
11  
12 520 transcription loop-mediated isothermal amplification (RT-LAMP), reverse transcription  
13  
14 521 polymerase chain reaction (RT-PCR), and their diagnostic results on MinION™ nanopore  
15  
16 522 sequencer for the detection of dengue virus serotypes. *Am J Microbiol Res.* 2015;3:118-24.  
17  
18  
19  
20 523 11. Hargreaves AD, Mulley JF. Snake venom gland cDNA sequencing using the Oxford nanopore  
21  
22 524 MinION portable DNA sequencer. *PeerJ.* 2015 Nov 24;3:e1441. doi: 10.7717/peerj.1441.  
23  
24 525 eCollection 2015.  
25  
26  
27  
28 526 12. Bolisetty MT, Rajadinakaran G, Graveley BR. Determining exon connectivity in complex  
29  
30 527 mRNAs by nanopore sequencing. *Genome Biol.* 2015;16:204.  
31  
32  
33 528 13. Cao MD, Ganesamoorthy D, Elliott A, Zhang H, Cooper M, Coin L. Streaming algorithms for  
34  
35 529 identification of pathogens and antibiotic resistance potential from real-time MinION  
36  
37 530 sequencing. *bioRxiv* 2015; doi: <http://dx.doi.org/10.1101/019356>.  
38  
39  
40  
41 531 14. Judge K, Harris SR, Reuter S, Parkhill J, Peacock SH. Early insights into the potential of the  
42  
43 532 Oxford Nanopore MinION for the detection of antimicrobial resistance genes. *J Antimicrob*  
44  
45 533 *Chemother* 2015; doi:10.1093/jac/dkv206.  
46  
47  
48 534 15. Wang J, Moore NE, Deng Y-M, Eccles DA, Hall RJ. MinION nanopore sequencing of an  
49  
50 535 influenza genome. *Front Microbiol.* 2015;6:766.  
51  
52  
53  
54 536 16. Kilianski A, Haas JL, Corriveau EJ, Liem AT, Willis KL, et al. Bacterial and viral identification  
55  
56 537 and differentiation by amplicon sequencing on the MinION nanopore sequencer. *GigaScience*  
57  
58 538 2015;4:12.  
59  
60  
61  
62  
63  
64  
65

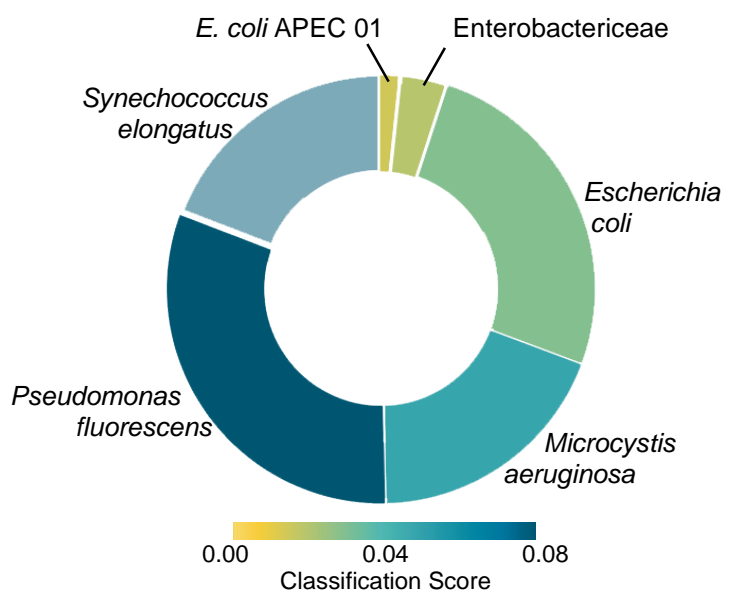
- 539 17. Greninger AL, Naccache SN, Federman S, Yu G, Mbala P, et al. Rapid metagenomic  
1 identification of viral pathogens in clinical samples by real-time nanopore sequencing  
2 540  
3  
4 541 analysis. *Genom Med.* 2015;7:99.  
5  
6  
7 542 18. Edwards A, Debbonaire AR, Sattler B, Mur LAJ, Hodson AJ. Extreme metagenomics using  
8  
9  
10 543 nanopore DNA sequencing: a field report from Svalbard, 78 °N. 2016. bioRxiv. doi:  
11  
12 544 <http://dx.doi.org/10.1101/073965>  
13  
14  
15 545 19. Brown BL, Watson M, Minot SS, Rivera MC, Franklin RB. Whole genome data from synthetic  
16  
17 546 metagenomes. *GigaScience Database.* 2017. <http://xxxxxx>.  
18  
19  
20 547 20. Watson M, Thomson M, Risse J, Talbot R, Santoyo-Lopez J, et al. poRe: an R package for the  
21  
22 548 visualization and analysis of nanopore sequencing data. *Bioinformatics.* 2015;31:114-5.  
23  
24  
25  
26 549 21. Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence data.  
27  
28 550 *Bioinformatics.* 2014; 30:3399–3401.  
29  
30  
31 551 22. Leggett RM, Heavens D, Caccamo M, et al.: NanoOK: Multi-reference alignment analysis of  
32  
33 552 nanopore sequencing data, quality, and error profiles. *Bioinformatics.* 2015;32:142-4.  
34  
35  
36 553 23. Meyer F, Paarmann D, D’Souza M, Olson R, Glass EM, et al. The metagenomics RAST server –  
37  
38 554 a public resource for the automatic phylogenetic and functional analysis of metagenomes.  
39  
40  
41 555 *BMC Bioinfo.* 2008;9:386.  
42  
43  
44 556 24. Juul S, Izquierdo F, Hurst A, Dai X, Wright A, Kulesha E, Pettett R, Turner D. What’s in my pot?  
45  
46 557 Real-time species identification on the MinION™. bioRxiv. 2015; doi:  
47  
48 558 <http://dx.doi.org/10.1101/030742>.  
49  
50  
51  
52 559 25. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact  
53  
54 560 alignments. *Genome Biol.* 2014;15:R46.  
55  
56  
57 561 26. Minot SS, Krumm N, Greenfield NB. One Codex: a sensitive and accurate data platform for  
58  
59 562 genomic microbial identification. bioRxiv. 2015. doi: <http://dx.doi.org/10.1101/027607>  
60  
61  
62  
63  
64  
65



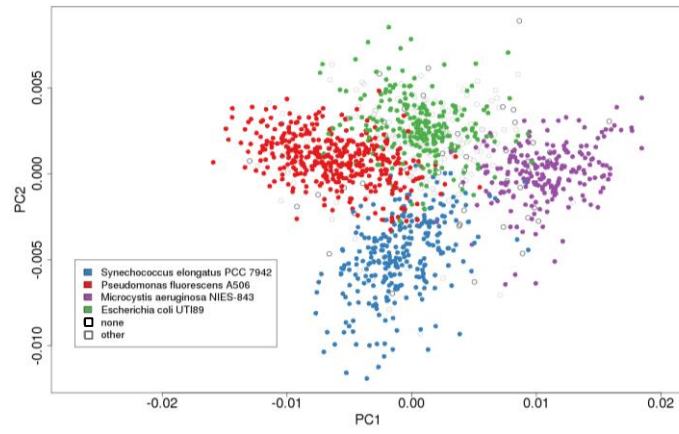
- 563 27. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST  
1 and PSI-BLAST: a new generation of protein database search programs. Nucl Acids Res.  
2 564  
3 and PSI-BLAST: a new generation of protein database search programs. Nucl Acids Res.  
4 565  
5 1997;25:3389.  
6
- 7 566 28. R Core Team. R: A language and environment for statistical computing. R Foundation for  
8  
9 Statistical Computing, Vienna, Austria. 2015;https://www.R-project.org/.  
10 567  
11
- 12 568 29. Park H-D, Sasaki Y, Maruyama T, Yanagisawa E, Hiraishi A, Kato K. Degradation of the  
13  
14 cyanobacterial hepatotoxin microcystin by a new bacterium isolated from a hypertrophic lake.  
15 569  
16 Environ Toxicol 2001;16:337-43.  
17 570  
18  
19
- 20 571 30. Jones MB, Highlander SK, Anderson EL, Li W, Dayrit M, et al. Library preparation methodology  
21  
22 can influence genomic and functional predictions in human microbiome research. Proc Nat  
23 572  
24 Acad Sci. 2015;45:14024-9.  
25 573  
26  
27
- 28 574 31. Wommack KE, Bhavsar J, Ravel J. Metagenomics: Read Length Matters. Appl Environ  
29  
30 Microbiol. 2008;74:1453.  
31  
32
- 33 576 32. Brown BL, LePrell RV, Franklin RB, Rivera MC, Cabral FC, et al. Metagenomic analysis of  
34  
35 planktonic microbial consortia from a non-tidal urban-impacted segment of James River.  
36 577  
37 Stand Genomic Sci. 2015;10:65.  
38 578  
39  
40
- 41 579 33. Magasin JD, Gerloff DL. Pooled assembly of marine metagenomic datasets: enriching annotation  
42  
43 through chimerism. Bioinformatics. 2015;31:311-317.  
44 580  
45
- 46 581 34. Freitas TAK, Li P-E, Scholz MB, Chain PSG. Accurate read-based metagenome characterization  
47  
48 using a hierarchical suite of unique signatures. Nucl Acids Res. 2015; doi:  
49 582  
50 10.1093/nar/gkv180.  
51 583  
52  
53
- 54 584 35. Zhang Q, Ye Y, Doak TG. Artificial functional difference between microbial communities caused  
55  
56 by length difference of sequencing reads. Biocomputing 2012;259-270. DOI:  
57 585  
58 http://dx.doi.org/10.1142/9789814366496\_0025.  
59 586  
60  
61  
62  
63  
64  
65

- 587 36. Frank JA, Pan Y, Tooming-Klunderud A, Eijsink VGH, McHardy AC, Nederbragt AJ, Pope PB.  
1  
2 588 Improved metagenome assemblies and taxonomic binning using long-read circular consensus  
3  
4 589 sequence data. *Sci Rep.* 2016;6:25373 doi:10.1038/srep25373.  
5  
6  
7 590 37. Sović I, Šikić M, Wilm A, Fenlon SN, Chen S, Nagarajan N. Fast and sensitive mapping of  
8  
9 591 nanopore sequencing reads with GraphMap. *Nat Comm.* 2016;11307.  
10  
11 592 doi:10.1038/ncomms11307 doi:10.1038/ncomms11307.  
12  
13  
14  
15 593 38. Chaisson MJ, Tesler G. Mapping single molecule sequencing reads using basic local alignment  
16  
17 594 with successive refinement (BLASR): application and theory. *BMC Bioinfo.* 2012;13:238  
18  
19 595 doi: 10.1186/1471-2105-13-238.  
20  
21  
22 596 39. Hugh R. Note: *Pseudomonas maltophilia* sp. nov., nom. rev. *Int J System Evol Microbiol.*  
23  
24 597 1981;31:195. doi: 10.1099/00207713-31-2-195  
25  
26  
27  
28 598 40. Binga EK, Lasken RS, Neufeld JD. Something from (almost) nothing: the impact of multiple  
29  
30 599 displacement amplification on microbial ecology. *The ISME Journal.* 2008;2:233-241.  
31  
32  
33 600 41. Lasken RS, Stockwell TB. (2007). Mechanism of chimera formation during the Multiple  
34  
35 601 Displacement Amplification reaction. *BMC Biotechnol* 7: 19. doi:10.1186/1472-6750-7-19  
36  
37  
38 602 42. Li C, Chng KR, Boey EJH, Ng AHQ, Wilm A, Nagarajan N. INC-Seq: accurate single molecule  
39  
40 603 reads using nanopore sequencing. *GigaScience.* 2016;5:34.  
41  
42  
43  
44 604 43. Benítez-Páez A, Portune KJ, Sanz Y. Species-level resolution of 16S rRNA gene amplicons  
45  
46 605 sequenced through the MinION™ portable nanopore sequencer. *GigaScience.* 2016;5:4.  
47  
48  
49 606 44. Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz M, McCombie WR. Oxford  
50  
51 607 Nanopore sequencing, hybrid error correction, and *de novo* assembly of a eukaryotic genome.  
52  
53 608 *Genome Res.* 2015b;25:1750-6.  
54  
55  
56  
57 609 45. Karlsson E, Lärkeryd A, Sjödin A, Forsman M, Stenberg P. Scaffolding of a bacterial genome  
58  
59 610 using MinION nanopore sequencing. *Scientif Rep* 2015;5:11996.  
60  
61  
62  
63  
64  
65

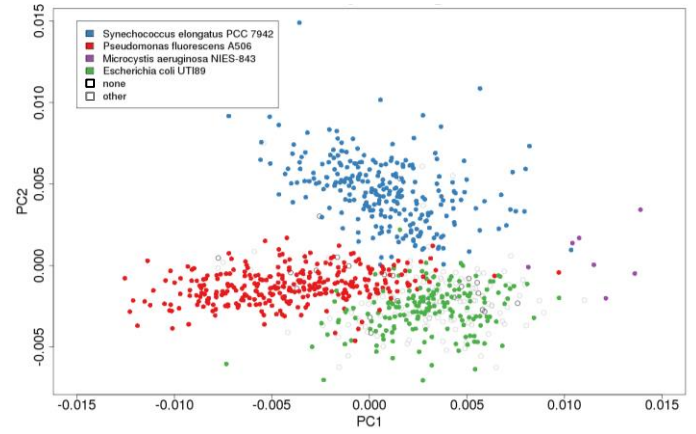
611 46. Stephens ZD, Lee SY, Faghri F, Campbell RH, Zhai C, et al. Big data: astronomical or  
1  
2 612 genomics? PLoS Biol. 2015;13:e1002195. doi:10.1371/journal.pbio.1002195.  
3  
4  
5 613 47. Garmendia C, Bernad A, Esteban JA, Blanco L, Salas M. The bacteriophage phi 29 DNA  
6  
7 614 polymerase, a proofreading enzyme. J Biol Chem. 1992;267: 2594-2599.  
8  
9  
10 615 48. Ip CLC, Loose M, Tyson JR, Cesare M, Brown BL, et al. MinION Analysis and Reference  
11  
12 Consortium: Phase 1 data release and analysis. F1000Research 2015;4:1075.  
13 616  
14  
15  
16 617 49. Kent WJ. BLAT—the BLAST-like alignment tool. Genome Res. 2002;12:656–64.  
17  
18  
19 618 50. Bartram AK, Lynch MD, Stearns JC, Moreno-Hagelsieb G, Neufeld JD. Generation of  
20  
21 619 multimillion-sequence 16S rRNA gene libraries from complex microbial communities by  
22  
23 620 assembling paired-end Illumina reads. Appl Environ Micro. 2011;77:3846–52.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

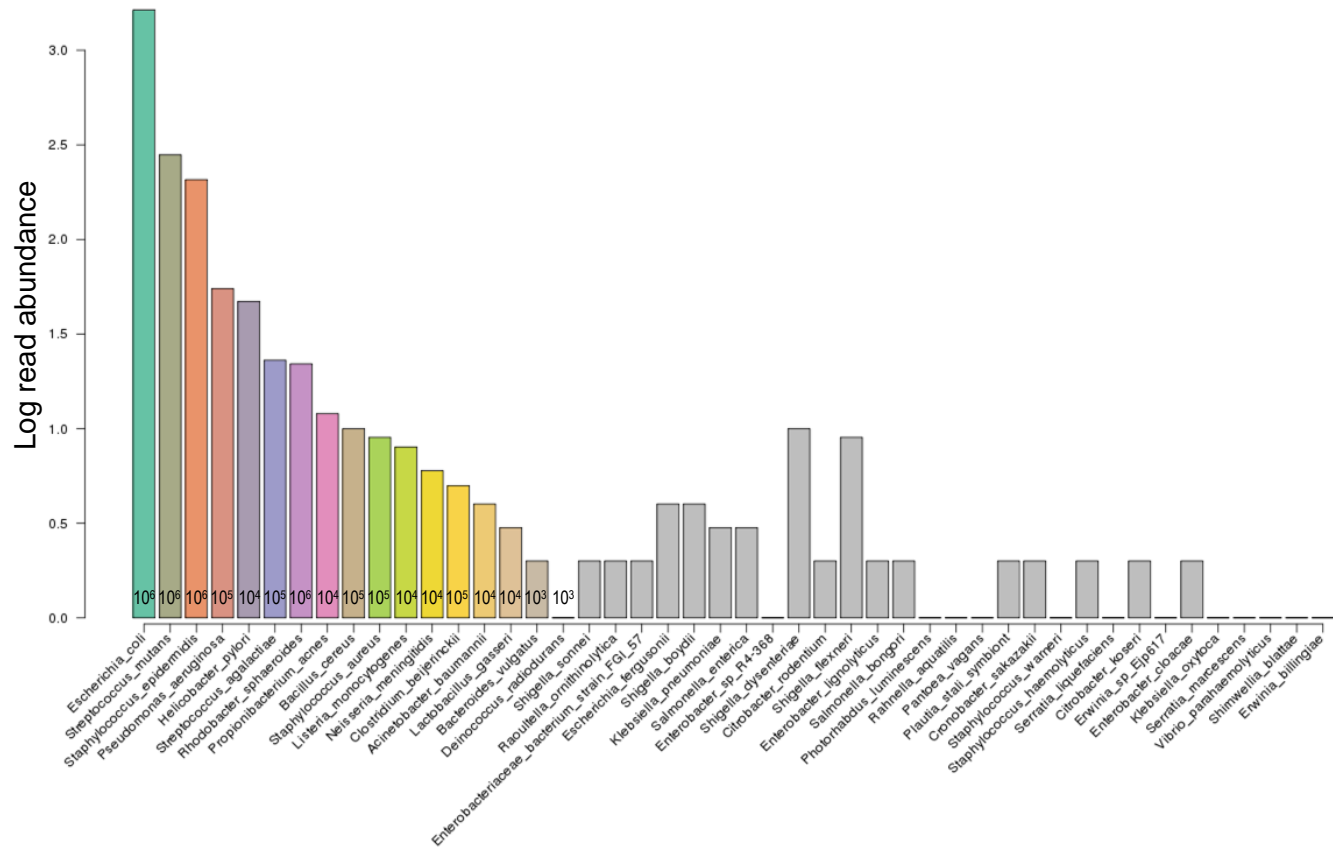


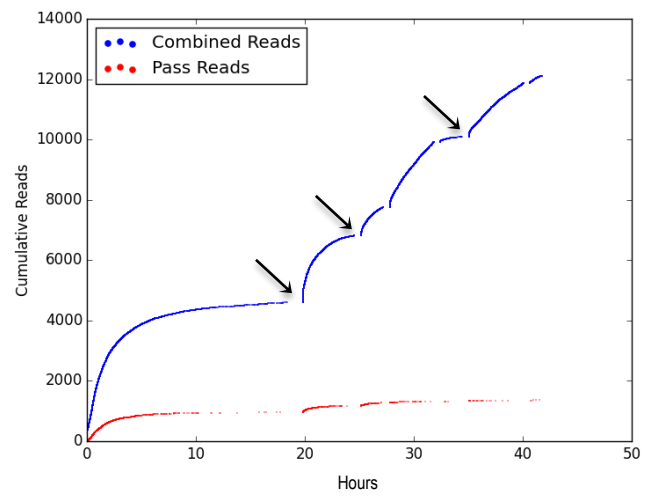
A



B









# VCU

Virginia Commonwealth University  
College of Humanities and Sciences  
Department of Biology  
Ecological Genetics Laboratory

1000 West Cary Street  
P.O. Box 842012  
Richmond, Virginia 23284-2012  
804-828-3265 voice  
804-828-0503 fax  
800-828-1120 TDD

6 January 2017

We respectfully submit the attached revised manuscript entitled, “MinION™ nanopore sequencing of environmental metagenomes: a synthetic approach,” for consideration for publication in *GigaScience*. Each and every comment and concern of the four reviewers has been addressed. As a result of the revision we have eliminated one set of results and added a different analysis, which in turn resulted in a change of authorship.

We have selected *GigaScience* as the venue to share our results with the world for two primary reasons. First, a number of other MinION papers are located in this journal, so our work will likely be seen by others who are considering use of this sequencing platform. Second, because a number of other metagenomic studies are published in this journal, we believe that our work will be noted by researchers who are likely to want to know how well MinION performs for metagenomic analysis.

Our manuscript is worthy of publication in *GigaScience* because it is well-conceived, considers in depth the consequences of long read sequencing of both simple and complex mixtures, and provides clear guidance for selecting analytical tools for long read sequence data analysis.

We take no exception to any issues relating to journal policies.

We have declared in the manuscripts our potential competing interests as follows: the contact author, Bonnie Brown, along with two other VCU authors (Maria Rivera and Rima Franklin) are members of the VCU MinION Access Programme. Mick Watson also is a member of this program. Sam Minot is not and has no other identifiable competing interests unless one considers that he is an employee of One Codex.

All authors have contributed to writing and have approved the manuscript for submission.

The content of this manuscript has not been published, nor has it been submitted for publication, elsewhere.

Sincerely,

A handwritten signature in black ink that reads 'Bonnie L. Brown'.

Bonnie L. Brown, Ph.D.

Professor and Associate Chair of Biology

Director, VCU Ecological Genetics Laboratory