1    Metagenomic characterization of a harmful algal bloom using nanopore sequencing

2

3    Peter W. Schafran[*], Victor Cai[†], Hsiao-Pei Yang[*], Fay-Wei Li[*,‡]

4    [*]Boyce Thompson Institute, Ithaca, NY 14853, USA

5    [†]Department of Biology, Duke University, Durham, NC 27708, USA

6    [‡]Plant Biology Section, Cornell University, Ithaca, NY 14853, USA

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27    Running Title: Nanopore HAB metagenome

28    Keywords: Nanopore, long-read sequencing, metagenomics, cyanobacteria, algal bloom,

29    Corresponding author:

30    Fay-Wei Li

31    Boyce Thompson Institute

32    533 Tower Rd.

33    Ithaca, NY 14853 USA

34    607-254-1244

35    fl329@cornell.edu

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53                                              ABSTRACT

54          Water bodies around the world are increasingly threatened by harmful algal blooms

55     (HABs) under current trends of rising water temperature and nutrient load. Metagenomic

56     characterization of HABs can be combined with water quality and environmental data to better

57     understand and predict the occurrence of toxic events. However, standard short-read

58     sequencing typically yields highly fragmented metagenomes, preventing direct connection of

59     genes to a single genome. Using Oxford Nanopore long-read sequencing, we were able to

60     obtain high quality metagenome-assembled genomes, and show that dominant organisms in a

61     HAB are readily identified, though different analyses disagreed on the identity of rare taxa.

62     Genes from diverse functional categories were found not only in the most dominant genera, but

63     also in several less common ones. Using simulated datasets, we show that the Flongle flowcell

64     may provide an option for HAB monitoring with less data, at the expense of failing to detect rarer

65     organisms and increasing fragmentation of the metagenome. Based on these results, we

66     believe that Nanopore sequencing provides a fast, portable, and affordable method for studying

67     HABs.

68

69                                             INTRODUCTION

70          Harmful algal blooms (HABs) pose serious threats to not only the aquatic biomes but

71     also human health. HABs occur when toxic or ecosystem-disrupting species of algae and

72     cyanobacteria rapidly increase in abundance in water bodies, leading to contamination of

73     drinking water supplies and closure of recreational waterways and fisheries. The presence of

74     HABs is associated with both acute toxicity, as well as chronic health issues such as cancer

75     (Etheridge 2010, Gorham et al. 2020). The frequency of HABs has increased with warming

76     water temperatures over the last few decades and may increase with additional change in

77     temperature, acidification, nutrient load, and oxygenation of fresh and marine aquatic

78     environments (Glibert et al. 2005). However, strong species and strain specific interactions

79    between the environment and aquatic microbial communities make prediction of HABs difficult

80    (Griffith and Gobler 2020), necessitating ongoing monitoring of species composition and

81    presence of toxins and taste-and-odor compounds, such as microcystins, geosmin, and 2-

82    methylisoborneol (Dietrich 2006). A survey of reservoir managers found that the control of algal

83    blooms was their highest concern (Schafran 2005). Considerable costs are incurred to control

84    HABs and their downstream effects, for example approximately $70 million spent over 10 years

85    in Waco, TX (Dunlap et al. 2015) and $80 million spent annually by the US Army Corps of

86    Engineers (Schafran 2005). Other estimates place costs of HABs at approximately $100 million

87    dollars annually in the US (Hoagland and Scatasta 2006).

88         Metabarcoding and metagenomic approaches to studying HABs are increasingly used

89    (Anderson 1995, Otten et al. 2016, Hennon and Dyhrman 2020, Hatfield et al. 2020), but as

90    generally applied have some limitations. Molecular probes and PCR primers are designed with

91    specific targets, and can introduce quantitative bias or may fail to amplify unknown members of

92    a microbial community (Krehenwinkel et al. 2017). The large (i.e. 23S/28S) and small (i.e.

93    16S/18S) ribosomal subunits and ribosomal internal transcribed spacers (ITS) have been most

94    commonly employed as a metabarcode (Lepère et al. 2000, Sherwood and Presting 2007, Stern

95    et al. 2012, Gamez 2018, Zhang et al. 2020), though primers targeting toxin and taste-and-odor

96    compound-producing genes have also been developed (Tsao et al. 2014, Suurnäkki et al. 2015,

97    Legrand et al. 2016). Metabarcoding and PCR-targeted approaches can only provide taxonomic

98    identity and presence/absence data for subjects of interest. Total metagenomic assembly allows

99    discovery of additional taxa and genes, though the typical 150-500 bp Illumina reads limit the

100   ability of assembly algorithms to link broad genomic regions. Library preparation and

101   sequencing typically rely on instrumentation, such as thermocyclers, qPCR machines, and DNA

102   sequencers, that are too large or unstable to function under field conditions. The time between

103   sample collection and data analysis is often longer than the 24-36 hours that can be necessary

104    to take action to prevent toxins and taste-and-odor compounds from infiltrating drinking water

105    systems.

106        The Oxford Nanopore MinION device offers a unique solution to overcome the above

107    limitations on HAB microbiome research. MinION is a third generation sequencing platform that

108    can produce an average of 15 Gbp of data consisting of long (> 10 kbp) reads, is portable, and

109    is powered by USB (Oxford Nanopore 2020). Compared to Illumina short-read sequencing

110    platforms with estimated error rates from 0.1-0.5% (Pfeiffer et al. 2018), the Nanopore error rate

111    is much higher at 5-15% (Rang et al. 2018), though this can be reduced to about 1% by

112    polishing the assembled contigs with read sequences (Vaser et al. 2017). The MinION's

113    portability has made it a popular choice for DNA sequencing where controlled laboratory

114    conditions are unavailable and/or real-time data are necessary, such as on Arctic/Antarctic

115    expeditions (Johnson et al. 2017, Gowers et al. 2019), on the International Space Station

116    (Castro-Wallace et al. 2017, Burton et al. 2020), in tropical rainforests (Menegon et al. 2017,

117    Pomerantz et al. 2018), and in clinical settings (Wongsurawat et al. 2019). Oxford Nanopore

118    recently released the inexpensive Flongle flow cell which allows for sequencing of up to 2 Gbp

119    at a cost of <20% that of the MinION flow cell (https://nanoporetech.com/products/ accessed 10

120    August 2020). This may be valuable in applications where increasing the number of sampling

121    events (e.g. long-term monitoring) is more valuable than deep sequencing depth.

122        In this study, we applied MinION to generate the first HAB metagenome on a water

123    sample from Cayuga Lake, New York. To assess consistency of taxonomic annotations among

124    the current tools on long-read DNA sequences, we evaluated out-of-the-box performance of

125    Centrifuge (Kim et al. 2016), Kaiju (Menzel et al. 2016), Kraken2 (Wood et al. 2019), and

126    BLAST+ (Camacho et al. 2009). We also created *in silico* datasets to mimic the Flongle flow cell

127    output and evaluate its consistency with MinION data.

128

129

130 <div align="center">MATERIALS AND METHODS</div>

131 **Sampling**

132       One surface water sample was collected from an ongoing HAB in Cayuga Lake at

133 Taughannock Falls State Park in Tompkins Co., NY (approximately 42.5470 N, 76.5984 W) on

134 July 15 2019. The water sample was split into 6 tubes (6ml each), and centrifuged at 21,000$g$

135 for 10min to pellet bacteria. For each tube, DNA was extracted using E.Z.N.A. Plant DNA kit

136 (Omega Bio-tek) following the manufacturer's protocol. The DNA was then pooled and

137 concentrated by AMPure XP beads to a total of 50µl for library prep.

138

139 **Sequencing and assembly**

140       The Nanopore sequencing library was prepared using the Ligation Sequencing kit (SQK-

141 LSK109) and sequenced on MinION R9 flowcell (FLO-MIN106D) for 60 hours. Signal files were

142 basecalled with Guppy v3.1.5 (Oxford Nanopore) with the high accuracy flip-flop mode.

143       Reads were adapter trimmed with Porechop (https://github.com/rrwick/Porechop) and

144 assembled using Flye v2.7 (Kolmorogov et al. 2019) in its metagenomic mode (--meta flag) with

145 an estimated genome size of 50 Mb (based on a previous assembly of the data). Assembled

146 contigs were polished with five iterations of Racon v1.4.3 (Vaser et al. 2017) followed by one

147 round of polishing with Medaka (https://nanoporetech.github.io/medaka). Assembly graph

148 structure was visualized with Bandage (Wick et al. 2015).

149

150 **Taxonomic assignment and gene annotation**

151       Polished contigs were taxonomically annotated using four tools in an out-of-the-box

152 fashion (i.e. without creation of custom search databases or extensive search parameter

153 optimization) in order to most closely replicate usage by end users. BLAST+ was used with the

154 NCBI *nt* database downloaded on 10 June 2020. Centrifuge was used with the NCBI nucleotide

155 non-redundant database last updated 3 March 2018 available from its website

156   (https://ccb.jhu.edu/software/centrifuge/). Kaiju was run with its *nr_euk* database, which is the

157   most inclusive and includes a subset of the NCBI *nr* database with sequences belonging to

158   archaea, bacteria, viruses, fungi, and microbial eukaryotes (https://github.com/bioinformatics-

159   centre/kaiju). Kraken2 was used with its option to create the NCBI *nt* database (kraken2-build --

160   download-library nt --db nt) which was constructed on 26 February 2020. Results were

161   interpreted based on either the best annotation for each contig selected by each program

162   (Centrifuge, Kaiju, Kraken2) or the result with the highest bitscore (BLAST+). NCBI taxonomy

163   database files (nodes.dmp and names.dmp) were used to parse taxonomic relationships of

164   annotations. Gene prediction was performed using MetaGeneMark (Besemer and Borodovsky

165   1999, Zhu et al. 2010), while ribosomal RNA genes were identified with barrnap 0.9

166   (https://github.com/tseemann/barrnap) and compared to SILVA 138 LSU/SSU Ref NR 99

167   databases (Quast et al. 2013, https://www.arb-silva.de/documentation/release-138/). rRNA

168   sequences were aligned with MAFFT v7.464 (Katoh and Standley 2013) and phylogenies

169   inferred with IQ-TREE v1.6.12 (Nguyen et al. 2015). Genes were annotated into clustered

170   orthologous groups (COGs) with eggNOG v5.0 (Huerta-Cepas et al. 2019). BUSCO 4.3.1 was

171   used to estimate completeness of gene annotation (Seppey et al. 2019).

172

173   **Flongle data simulation**

174   To mimic output by the Nanopore Flongle flowcell, data were subset to approximately 2

175   Gbp by pseudo-randomly selecting reads using shuf (GNU coreutils 8.21). Ten subsets were

176   generated and checked to ensure a similar length distribution to the complete dataset. Each

177   subset was assembled, polished, and annotated with BLASTN as described above.

178

179   **Data Availability**

180   Supplemental files available through Figshare

181   (https://figshare.com/projects/Nanopore_HAB_metagenome/92567) include the polished

182 metagenome (File S1, Predicted genes and their locations within the metagenome (File S2),

183 amino acid sequences (File S3), and taxonomic annotations for each contig at each taxonomic

184 level (File S4).

185               RESULTS

186   Following adapter trimming, 3.74 million reads (9.35 Gbp) were retained with an N50

187 length of 3651 bp and a mean length of 2498 bp. Maximum read length was 110 kbp. Assembly

188 and contig polishing produced a metagenome size of 48.9 Mbp consisting of 2183 contigs with

189 an N50 of 90 kbp (Table 1; File S1).

190

191 **Taxonomic Annotation**

192   Taxonomic annotation varied considerably among the classification tools. The number of

193 contigs that could not be classified or could only be classified at the root of the taxonomy ranged

194 from 255 (Centrifuge) to 500 (Kaiju). Because the unclassified contigs tended to be short, these

195 represent only 1.5% (Centrifuge, Kraken) to 15% (Kaiju) of the total length of the assembly

196 (Figure S1). At the superkingdom level, the majority of classified contigs were identified as

197 Bacteria (range 81-90%), while Eukaryota were represented by 10-18% of contigs. Archaea and

198 viruses made up a very small proportion of the results, representing less than 0.5% of classified

199 contigs. At the genus level, the number of taxa identified ranged from 263 (Kaiju) to 309

200 (Centrifuge). Only 53 genera were shared among all classification tools, representing between

201 43% to 77% of the total assembly, and of these only 37 genera individually represent >0.5% of

202 the contigs (Figure 1). In each case, over 50% of contigs were identified as *Anabaena* or

203 *Dolichospermum* (here combined due to recent taxonomic revisions between these genera;

204 Wacklin et al. 2009). Taxa shared similar proportions among tools at other taxonomic levels

205 (Figure S2, File S4). For individual contigs, 27% were identified to the same genus by all tools,

206 while for an additional 17% of contigs three of the tools agreed on one genus. For 20% of

207 contigs, there was no agreement between any classification tools (including failure to classify at

208    genus level). These ambiguously annotated contigs on average tended to be longer than the

209    average length of other contigs, with an N50 length of 152 kbp vs 90 kbp, respectively. The

210    number of contigs vs. total assembly size per genus showed a positive linear relationship (p <

211    $2e^{-16}$, Kendall's $\tau$ = 0.54).

212        Of the largest, high coverage cluster within the assembly graph—presumably

213    representing *Anabaena/Dolichospermum* genome(s)—526 contigs form a relatively 'knotty'

214    structure with many edges generally converging on two nodes (Figure 2). About 72% of contigs

215    in this cluster were annotated as *Anabaena* or *Dolichospermum*. Coloration of edges by

216    coverage suggests two levels, supported by Hartigan's dip test of unimodality (p < $2.2e^{-16}$). This

217    cluster may represent two genomes from closely related species, but attempts to isolate and

218    reassemble them into more contiguous genomes were unsuccessful. There are 47 long contigs

219    (>100 kbp) that are disconnected to the *Anabaena/Dolichospermum* cluster, which tend to have

220    a lower coverage, except one 648 kbp contig annotated as *Anabaena/Dolichospermum* above

221    50X coverage. Three circular contigs with coverage from 390-628 X and lengths from 86-95 kbp

222    were also annotated as *Anabaena/Dolichospermum*, and might be their plasmids.

223

**Gene Annotation**

225        Gene annotation identified 66,490 genes with an N50 length of 756 bp (Files S2, S3).

226    BUSCO analysis with the Bacteria v10 dataset found 99% complete genes in our metagenome

227    (including 88% in the *Anabaena/Dolichospermum* cluster), suggesting this annotation is fairly

228    robust. Genes were found representing every COG category (Figure 3). After genes with

229    unknown function, the largest COG represented genes involved with amino acid transport and

230    metabolism. Several COGs had very few representative genes, including: 1) RNA processing

231    and modification; 2) chromatin structure and dynamics; 3) extracellular structures; 4) nuclear

232    structure; 5) cytoskeleton. While the most abundant genera in the assembly also had the

233   highest number of genes, several genera that were relatively rare in the assembly showed a

234   high number of genes distributed across the predominant COG categories (Figure 3).

235       Barrnap identified thirty-seven 16S and thirty-two 23S sequences with average lengths

236   of 1185 bp and 2469 bp, respectively, though several fragmented sequences needed manual

237   correction. Twenty-three 16S and fifteen 23S sequences matched entries in SILVA databases

238   with greater than 95% identity and 95% coverage of the query sequence length. Phylogenies

239   showed the most common rRNA sequences are nested in *Anabaena* (Figures S3, S4). Pairs of

240   16S and 23S from the same contig supported the presence of *Pseudanabaena, Candidatus*

241   Nanopelagicus, *Candidatus* Fonsibacter, *Cercopagis pengoi,* and an unknown taxon sister to

242   *Curvibacter, Rhodoferax,* and *Polaromonas.* In other cases 16S/23S pairs conflicted in their

243   position to reference taxa. Several 16S/23S pairs nested within Eukaryota but with low similarity

244   to all SILVA references may represent organellar rRNA sequences. However, these sequences

245   also failed to match any sequence in NCBI *nt* with identity and alignment length above 90%.

246

247   **Flongle Data Simulations**

248       The 10 data subsets created to mimic output of the Flongle flowcell resulted in

249   assemblies with an average size of 17 Mbp (standard deviation = 0.4 Mbp) composed of an

250   average of 965 contigs (standard deviation = 46). While each subset represented 21% of the

251   total data, the subset assemblies averaged 35% of the size of the total assembly. Similarly, the

252   number of genes and total size of the genes both averaged 35% as large as the total assembly.

253       Taxonomic annotation by BLAST was largely consistent at the superkingdom level,

254   though a small percentage of the subset assemblies was identified as Archaea, which was not

255   present in the total assembly (Figure 4). For other superkingdoms, the total assembly falls within

256   the kernel density estimate of the subset results. At the genus level, larger discrepancies

257   appeared (Figure 4), with similar trends apparent at other taxonomic levels (Figure S5).

258   *Anabaena/Dolichospermum* was more highly represented in the subset assemblies than the

259      total assembly, while most other genera showed the opposite trend. The range of subset results

260      was large for some genera, such as *Sphingomonas* with a range of 15 percentage points,

261      *Sphingopyxis* with 11 percentage points, and unannotated contigs with 11 percentage points.

262

263                                  DISCUSSION

264            Our exploratory analysis using Nanopore DNA sequencing to characterize a harmful

265      algal bloom highlights some potential uses for this technology, as well as some areas where

266      caution should be used in interpreting the data. The main drawback of long-read DNA

267      sequencing technology is a relatively high error rate compared to short-read sequencers. The

268      error rate of raw sequences for the Nanopore flowcell and basecaller we used is estimated at

269      roughly 10% (Wick et al. 2019), too high to be useful for k-mer based analysis or taxonomic

270      annotation (data not shown). The development of Nanopore-focused assemblers that take into

271      account the difference in coverage typical of metagenomic data, as well as polishing tools that

272      remove errors through consensus calling and neural networks trained on specific models of

273      Nanopore errors, provides methods to create highly accurate metagenomes without short-read

274      sequences. These methods offer an advantage over short-read metagenomes by greatly

275      increasing the contiguity of assembly with a small trade-off in accuracy. Our assembly returned

276      many large contigs (N50 = 90 kbp), including one 2.3 Mbp contig identified as

277      Hyphomonadacae bacterium that is about two-thirds the length of its most closely related

278      bacterial genome in NCBI's Genome database. We expected a more contiguous genome of the

279      predominant cyanobacteria in our sample (*Anabaena* and/or *Dolichospermum*), but those

280      contigs formed a few highly connected and structurally unresolved clusters. This can occur

281      when a genome contains long repeat-rich regions or when multiple organisms share portions of

282      their genomes that are similar. The structure of the assembly graph combined with significant

283      bimodal distribution in coverage suggests that at least two highly similar organisms may be

284      present in the *Anabaena/Dolichospermum* cluster, resulting in several shared edges connecting

285    otherwise distinct clusters of contigs. However, our attempts to isolate and reassemble the raw

286    data contributing to the *Anabaena/Dolichospermum* cluster did not improve the assembly.

287        BLAST, Centrifuge, Kaiju, and Kraken2 all were able to provide taxonomic annotations

288    for a majority of contigs. Kaiju failed to annotate longer contigs than the others, possibly

289    because it uses a smaller protein search database derived from NCBI *nr*, as opposed to the

290    nucleotide-based *nt* used by the other tools. The low level of shared genus richness—where

291    only 27% of contigs were placed in the same genus by every tool and roughly 20% of the

292    genera identified by each tool are shared among all four—shows that caution needs to be

293    exercised in determining community composition. While the more abundant organisms in this

294    community may be identified with some certainty, contigs receiving mixed taxonomic

295    annotations may need to be manually inspected. We processed BLAST results using the match

296    between each contig and the reference database with the highest score to assign an identity to

297    the contig. Individual matches tended to be short, with an average of 537 bp, so it may be more

298    informative to examine all matches across different parts of the contig and use taxonomic

299    distance (least common ancestor) between best hits to refine the overall identity, similar to how

300    Centrifuge, Kaiju, and Kraken2 operate. Kraken2 provides an option to filter results by a user-

301    defined confidence score, requiring a given fraction of a sequence to match a taxonomic level to

302    receive an annotation. However, even at a low confidence level of 0.1, Kraken2 failed to provide

303    many meaningful results from our data with 68% of contigs unclassified and of those classified,

304    many were placed at the order level and above, suggesting that the application of this approach

305    for metagenomics is limited.

306        Any classification is only as good as the reference database being searched. While

307    NCBI *nt/nr* are the traditional databases for classification software, inclusion of the NCBI Whole

308    Genome Shotgun (WGS) database can greatly improve results (Martí and Garay 2019). The

309    short matches between our metagenome and *nt/nr* databases highlight the need to incorporate

310    more genomic data for comparison against the increasing number of long-read metagenomic

311    surveys. The software we tested all provide ways to create custom search databases, but this

312    can require technical knowledge and resources many end users do not possess (e.g. our server

313    with 768 GB RAM was memory-limited when trying to create a Centrifuge database from the

314    current *nt* database). As a result, software that are not provided with regularly updated search

315    databases, such as the most recent Centrifuge *nt* release from April 2018, quickly become

316    obsolete as new genomic data are regularly added to the source databases. Additionally,

317    updates in nomenclature, such as the transfer of some *Anabaena* into *Dolichospermum,* and

318    persistence of polyphyly in bacterial phylogenies will confound attempts to use a least common

319    ancestor approach to reconcile multiple matches (Wacklin et al. 2009, Cirés and Ballot 2016).

320    This may partially explain why the proportion of contigs placed into *Anabaena* vs.

321    *Dolichospermum* varies between BLAST/Kraken2 vs. Centrifuge/Kaiju, but in combination the

322    genera are roughly equal. An independent sample from the same HAB event was

323    microscopically identified to *Dolichospermum* (Bloom code 19-3435-B2, Community Science

324    Institute, Ithaca, NY), supporting our decision to combine *Anabaena* and *Dolichospermum*

325    results.

326        Our gene annotation and classification showed a broad diversity of genes obtained not

327    only from the most abundant organisms, but also from many with lower abundance in our

328    sample. This demonstrates one of the advantages of long-read sequencing in that much longer

329    contigs allow connecting more genes to a single taxon, in addition to being able to taxonomically

330    place previously unknown genes or genes without a known function. Efforts to incorporate

331    genomic and proteomic data into HAB prediction models will be aided by this technologic

332    development (Hennon and Dyhrman 2020). While we did not identify any in our sample, genes

333    related to production of toxin and taste-and-odor compounds, such as geosmin synthase, 2-

334    methylisoborneol (MIB) synthase, the anatoxin-a synthetase gene cluster, and the microcystin

335    synthetase gene cluster, could be easily gleaned from this type of metagenome. This is

336    consistent with failure to detect microcystin in this HAB at the drinking water threshold (<3 µg/L)

337     by the Community Science Institute (Ithaca, NY). The ability of the Nanopore platform to also

338     sequence RNA suggests that monitoring the expression of these genes could provide a real-

339     time, field-based method to track the production of these dangerous compounds.

340         We created ten datasets *in silico* to estimate the potential of Oxford Nanopore's Flongle

341     flow cell to assemble and characterize a metagenome as well as the MinION flow cell that

342     produced 5 times more data. The average Flongle assembly size was disproportionately large

343     compared to the reduced size of the datasets, suggesting that some of the MinION data could

344     be considered excess. However, the distribution of contig length shows the Flongle assemblies

345     varied tremendously, with the length of the longest contig ranging from 2.5 to 0.6 Mbp and N50

346     ranging from 60 to 78 Kbp. Since the length distribution of the reads in each dataset was similar,

347     this suggests the difference in assemblies is due to random assortment of reads derived from

348     different organisms. While the proportion of the Flongle assemblies annotated to each

349     taxonomic group was similar to the MinION assembly at the superkingdom level, increasing

350     levels of taxonomic resolution showed the most prevalent organisms were overrepresented in

351     the Flongle results, while the least prevalent organisms tended to be underrepresented. This is

352     consistent with the methodology of the assembly and polishing software, where contigs with

353     very low coverage are removed from the final assembly. Based on our results the Flongle may

354     be useful for capturing large genomic fractions of the predominant members of an HAB, though

355     it is more likely to miss minor taxa in the community.

356

357     **Conclusions**

358         The profound effects of harmful algal blooms (HABs) on human health, economic

359     activity, and ecological function compounded with the expectation that changes to the

360     environment will alter the frequency and intensity of HABs necessitates a better understanding

361     of the patterns and processes driving these events. Genomics provide complementary evidence

362     to traditional data focused on taxonomic identification, quantification, and water quality

363    measurements. Comparing four tools used for classifying DNA sequences, we found that

364    assigning taxonomy with high confidence in our Nanopore-derived HAB metagenome was often

365    difficult, though as more complete genomes are added to reference databases, we expect this

366    to be resolved. We recovered a diversity of genes from many functional groups, even from

367    organisms present at relatively low abundance in the metagenome. Our *in silico* datasets

368    mimicking the Flongle flow cell for the Nanopore MinION and GridION sequencers suggests this

369    may be a valuable tool for water quality managers who need to repeatedly monitor for HABs,

370    with the tradeoffs being decreased detection of low abundance organisms and greater

371    fragmentation of the metagenome. Given the portability, affordability, and short time from

372    sample-to-sequence of the Nanopore platform, we believe this form of long-read sequencing

373    shows potential for the study and monitoring of algal blooms.

374

375                              ACKNOWLEDGEMENTS

383

384                              CONFLICTS OF INTEREST

385    The authors declare no conflicts of interest.

386

387

388

389 **Table 1**. Summary statistics of assemblies and predicted genes. Flongle results represent mean

390 (standard deviation) of all 10 subsets.

| Assembly | Total Length (Mbp) | Number | Mean Length (bp) | Longest (Mbp) | Shortest (bp) | N50 | N70 | N90 |
|---|---|---|---|---|---|---|---|---|
| MinION | 48.9 | 2183 | 22408 | 2.30 | 80 | 90015 | 42028 | 12372 |
| Flongle | 17.1 (0.441) | 965 (46.4) | 17773 (754) | 1.31 (0.625) | 67 (28) | 64864 (4341) | 37296 (1432) | 10570 (1369) |

| Predicted Genes | Total Length (Mbp) | Number | Mean Length (bp) | Longest (bp) | Shortest (bp) | N50 | N70 | N90 |
|---|---|---|---|---|---|---|---|---|
| MinION | 12.3 | 66490 | 184.8 | 4721 | 18 | 252 | 166 | 90 |
| Flongle | 4.28 (0.103) | 23268 (720) | 184.2 (1.92) | 3379 (572) | 18 (0) | 255.3 (2.75) | 164.1 (2.02) | 89.1 (1.10) |

391

392

393

394

395

**Figure 1**. Proportion of contigs classified to superkingdoms (A) and genera (B) by each classification tool. Only genera identified by all classification tools that represent more than 0.5% of assembly shown. Individual genera present at <0.5% were combined into 'Other'. Genera grouped by phylum.

**Figure 2**. Assembly graph colored by coverage. Five largest clusters annotated with taxon at lowest level that provides majority agreement among all contigs. Contig lengths to scale.

**Figure 3**. 2D histogram of log-transformed COG annotations for genera with more than 100 annotated genes. Color legend represents log10-transformed number of genes. Histograms represent number of genes in respective row/column. X-axis categories described in left-hand legend.

**Figure 4**. Comparison of assembly sizes for superkingdoms and genera between total assembly (dots) and Flongle-mimic subsets (violin plots) as annotated by BLAST. Only genera representing >1% of assembly shown.

**Figure S1**. Histograms of unclassified contigs by length. Blue = Blast, Red = Centrifuge, Green = Kaiju, Yellow = Kraken2, White = all contigs.

**Figure S2**. Proportion of contigs classified to phylum (A) and family (B) by each classification tool. Only taxonomic groups identified by all classification tools shown.

A

Heterosigma akashiwo GQ222227.36970.38445
Gelidium vagum KC875854.23195.24567
16S_rRNA__contig_1427_segment0_46126–47826
16S_rRNA__contig_1427_segment0_5–1710
16S_rRNA__contig_1516_segment1_1211–2155
16S_rRNA__contig_1333_segment0_1–1423
16S_rRNA__contig_1610_segment0_23261–24633
16S_rRNA__contig_1591_segment0_53680–54650
uncultured eukaryote KJ925287.1.1638
uncultured phyllopharyngid ciliate AY821935.1.1508
16S_rRNA__contig_1065_segment0_1947–4432
uncultured eukaryote KY355500.1.1609
Acineta flava HM140400.1.1706
uncultured marine picoeukaryote FR874559.1.1705
uncultured marine picoeukaryote FR874828.1.1704
16S_RNA__contig_1763_segment0_2133–3895
uncultured ciliate HQ219368.1.2292
uncultured microeukaryote JN705509.1.1368
metagenome FPLS01057222.25.1303
uncultured alveolate DQ244029.1.1753
Pseudocyrtolophosis alpestris EU264564.1.1760
uncultured marine eukaryote HM581713.1.1778
uncultured eukaryote KJ759177.1.1775
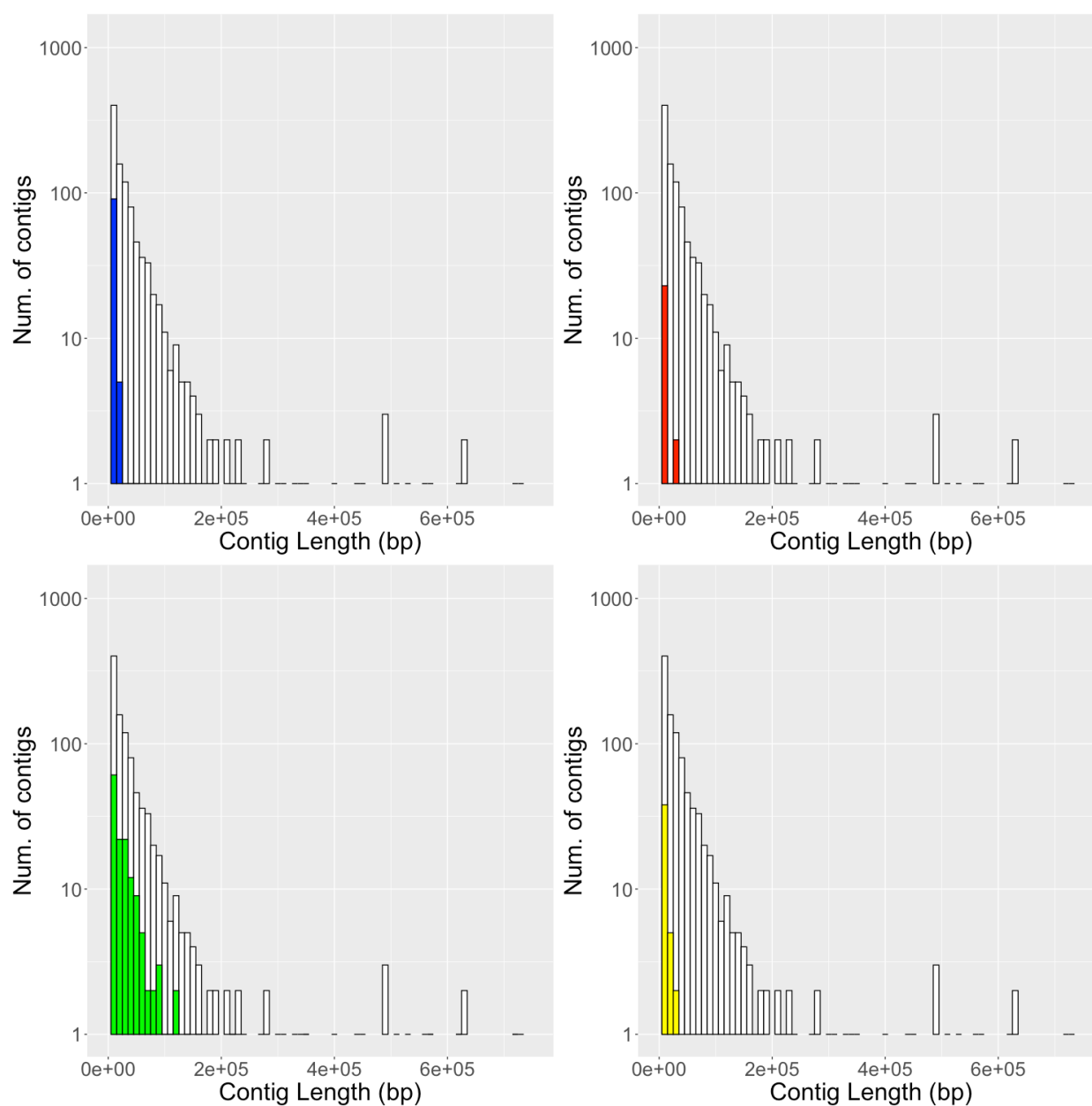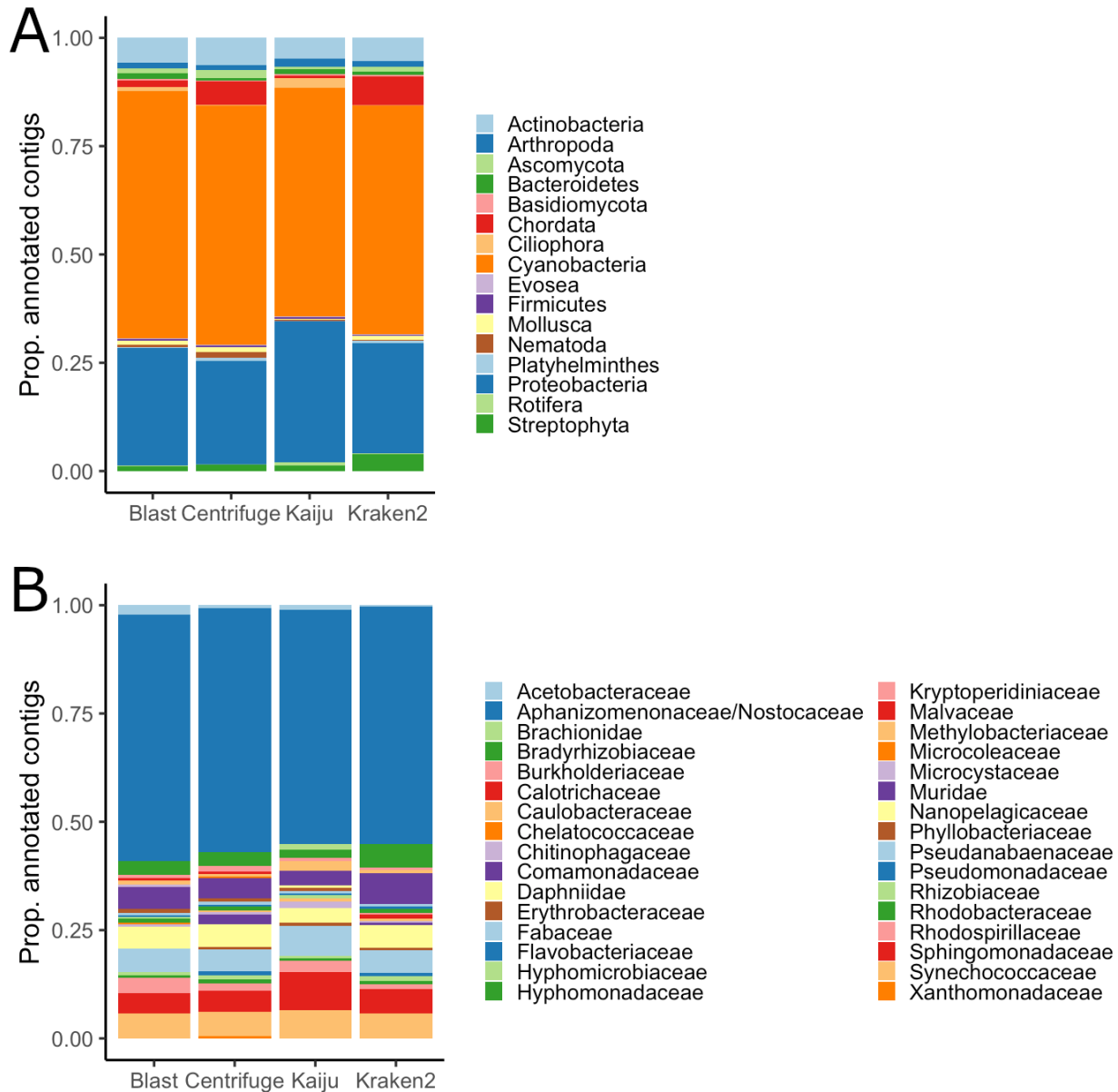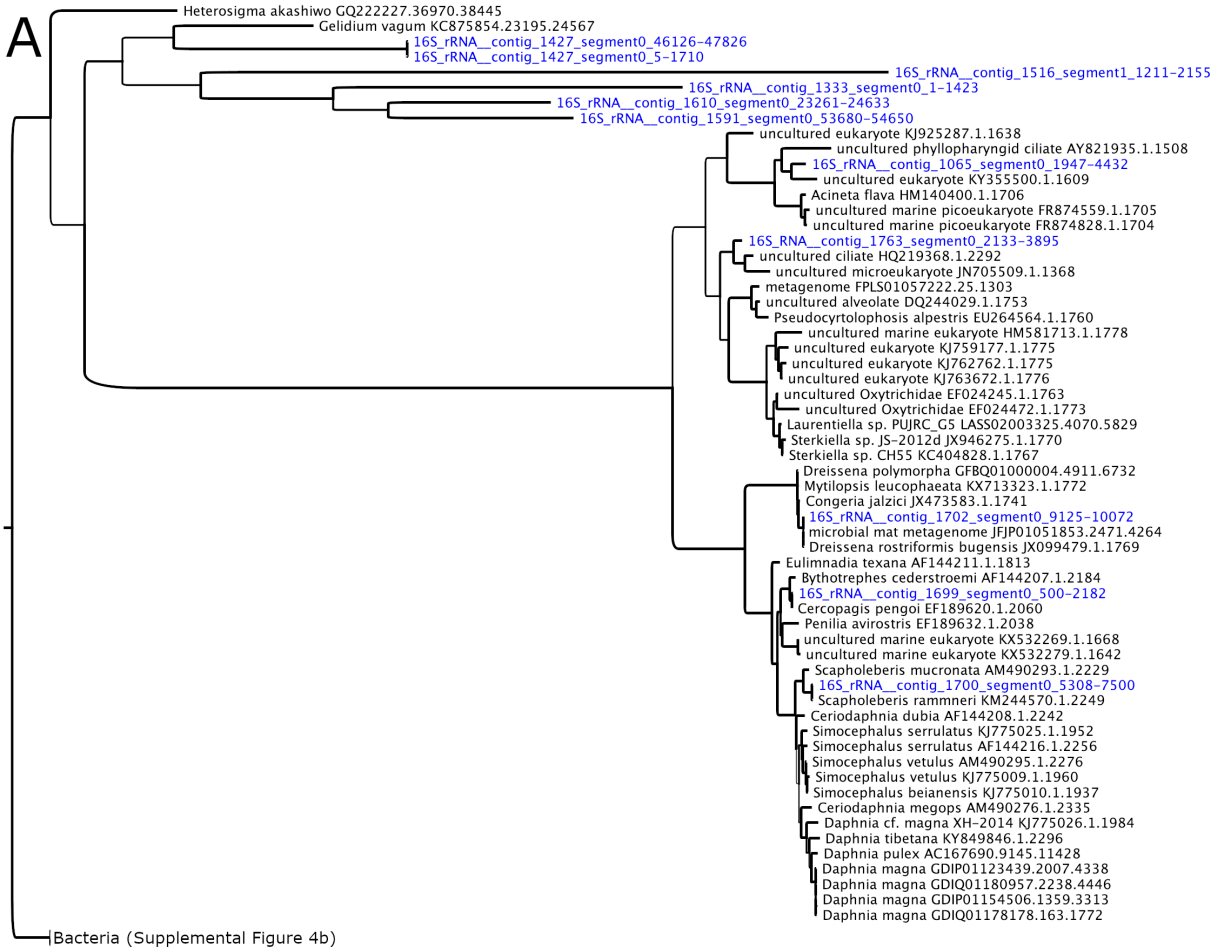uncultured eukaryote KJ762762.1.1775
uncultured eukaryote KJ763672.1.1776
uncultured Oxytrichidae EF024245.1.1763
uncultured Oxytrichidae EF024472.1.1773
Laurentiella sp. PUJRC_G5 LASS02003325.4070.5829
Sterkiella sp. JS–2012d JX946275.1.1770
Sterkiella sp. CH55 KC404828.1.1767
Dreissena polymorpha GFBQ01000004.4911.6732
Mytilopsis leucophaeata KX713323.1.1772
Congeria jalzici JX473583.1.1741
16S_rRNA__contig_1702_segment0_9125–10072
microbial mat metagenome JFJP01051853.2471.4264
Dreissena rostriformis bugensis JX099479.1.1769
Eulimnadia texana AF144211.1.1813
Bythotrephes cederstroemi AF144207.1.2184
16S_rRNA__contig_1699_segment0_500–2182
Cercopagis pengoi EF189620.1.2060
Penilia avirostris EF189632.1.2038
uncultured marine eukaryote KX532269.1.1668
uncultured marine eukaryote KX532279.1.1642
Scapholeberis mucronata AM490293.1.2229
16S_rRNA__contig_1700_segment0_5308–7500
Scapholeberis rammneri KM244570.1.2249
Ceriodaphnia dubia AF144208.1.2242
Simocephalus serrulatus KJ775025.1.1952
Simocephalus serrulatus AF144216.1.2256
Simocephalus vetulus AM490295.1.2276
Simocephalus vetulus KJ775009.1.1960
Simocephalus beianensis KJ775010.1.1937
Ceriodaphnia megops AM490276.1.2335
Daphnia cf. magna XH–2014 KJ775026.1.1984
Daphnia tibetana KY849846.1.2296
Daphnia pulex AC167690.9145.11428
Daphnia magna GDIP01123439.2007.4338
Daphnia magna GDIQ01180957.2238.4446
Daphnia magna GDIP01154506.1359.3313
Daphnia magna GDIQ01178178.163.1772

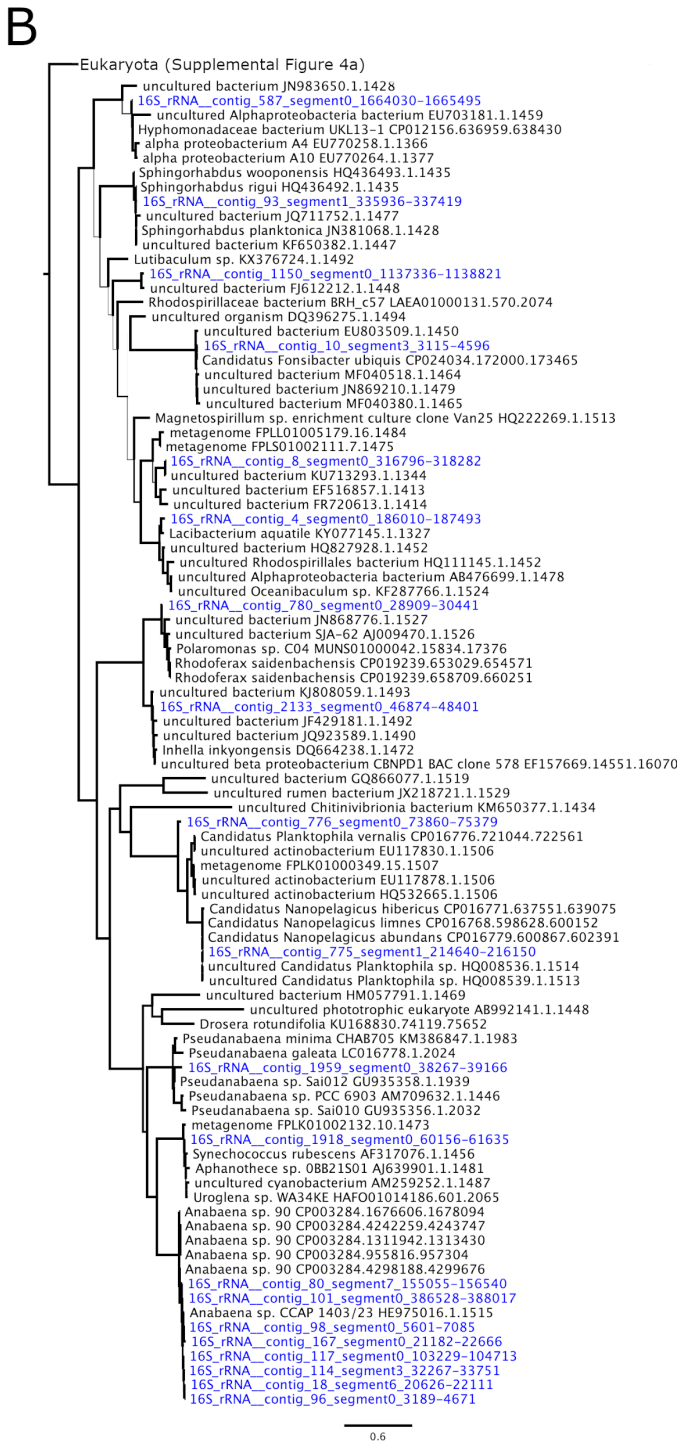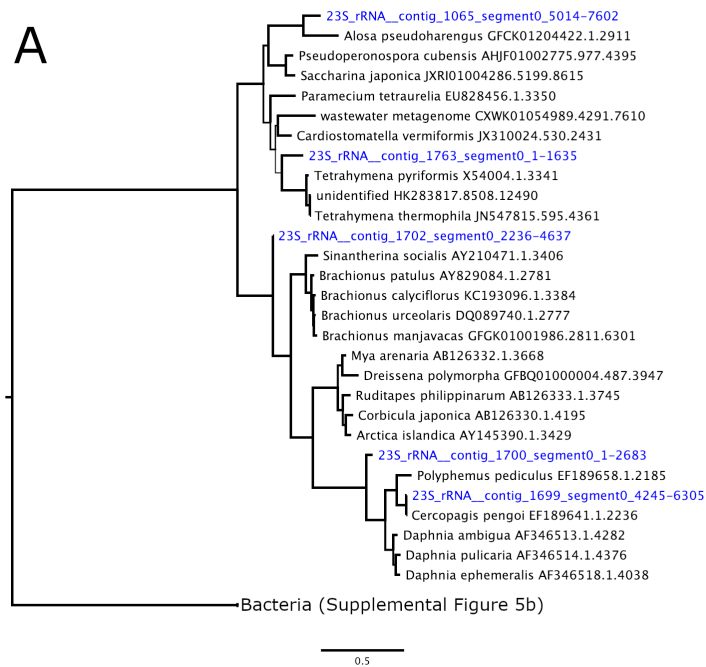Bacteria (Supplemental Figure 4b)

434

0.6

435

B

**Figure S3**. 16S rRNA phylogeny of metagenome-recovered sequences (blue) and SILVA

references (black). Tree rooted along the branch separating Eukaryota (A) and Bacteria (B).

Branch thickness is proportional to bootstrap support with scale bar width = 100.
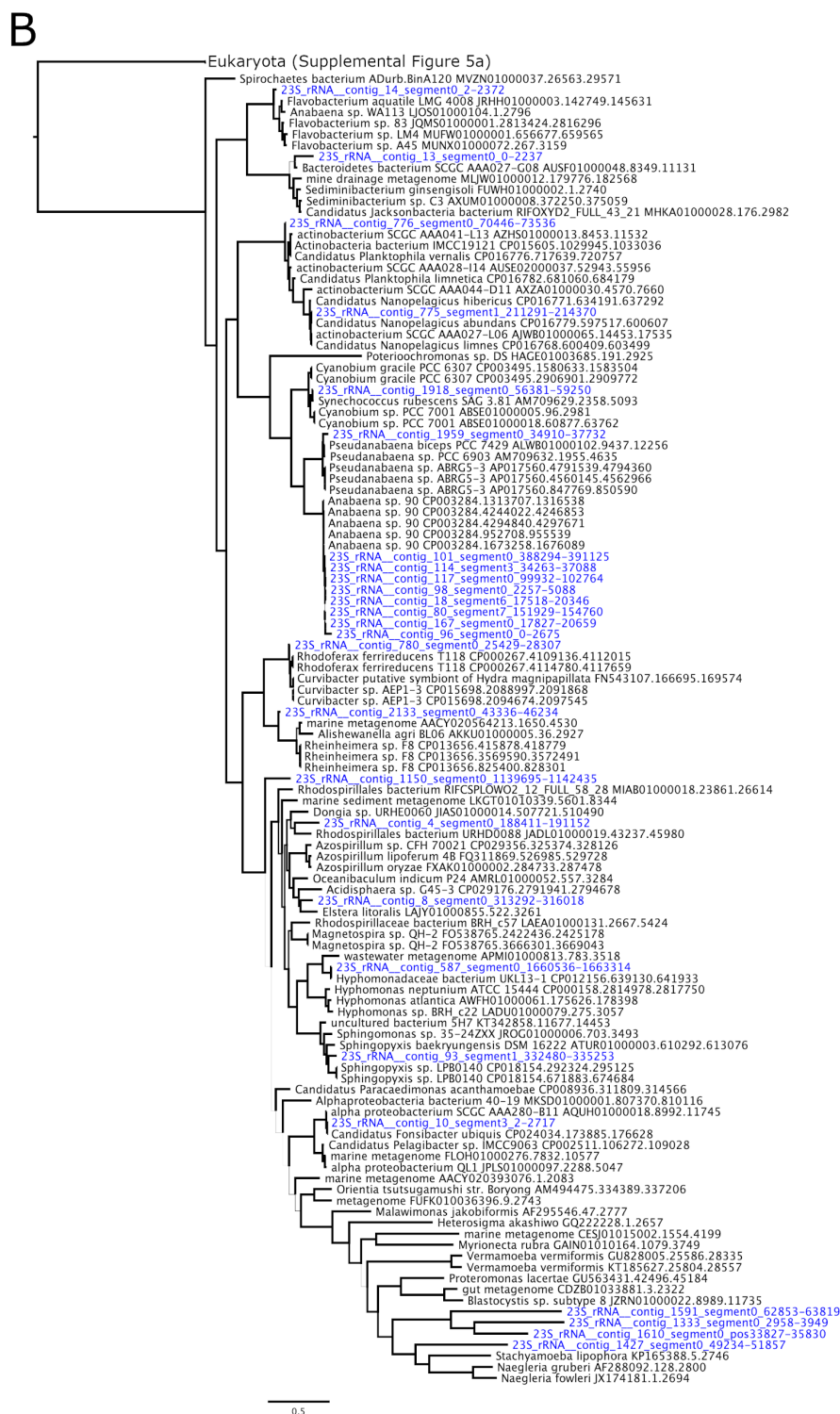
A

23S_rRNA__contig_1065_segment0_5014–7602
Alosa pseudoharengus GFCK01204422.1.2911
Pseudoperonospora cubensis AHJF01002775.977.4395
Saccharina japonica JXRI01004286.5199.8615
Paramecium tetraurelia EU828456.1.3350
wastewater metagenome CXWK01054989.4291.7610
Cardiostomatella vermiformis JX310024.530.2431
23S_rRNA__contig_1763_segment0_1–1635
Tetrahymena pyriformis X54004.1.3341
unidentified HK283817.8508.12490
Tetrahymena thermophila JN547815.595.4361
23S_rRNA__contig_1702_segment0_2236–4637
Sinantherina socialis AY210471.1.3406
Brachionus patulus AY829084.1.2781
Brachionus calyciflorus KC193096.1.3384
Brachionus urceolaris DQ089740.1.2777
Brachionus manjavacas GFGK01001986.2811.6301
Mya arenaria AB126332.1.3668
Dreissena polymorpha GFBQ01000004.487.3947
Ruditapes philippinarum AB126333.1.3745
Corbicula japonica AB126330.1.4195
Arctica islandica AY145390.1.3429
23S_rRNA__contig_1700_segment0_1–2683
Polyphemus pediculus EF189658.1.2185
23S_rRNA__contig_1699_segment0_4245–6305
Cercopagis pengoi EF189641.1.2236
Daphnia ambigua AF346513.1.4282
Daphnia pulicaria AF346514.1.4376
Daphnia ephemeralis AF346518.1.4038
Bacteria (Supplemental Figure 5b)

0.5

441

442

**Figure S4**. 23S rRNA phylogeny of metagenome-recovered sequences (blue) and SILVA

references (black). Tree rooted along the branch separating Eukaryota (A) and Bacteria (B).

Branch thickness is proportional to bootstrap support with scale bar width = 100.
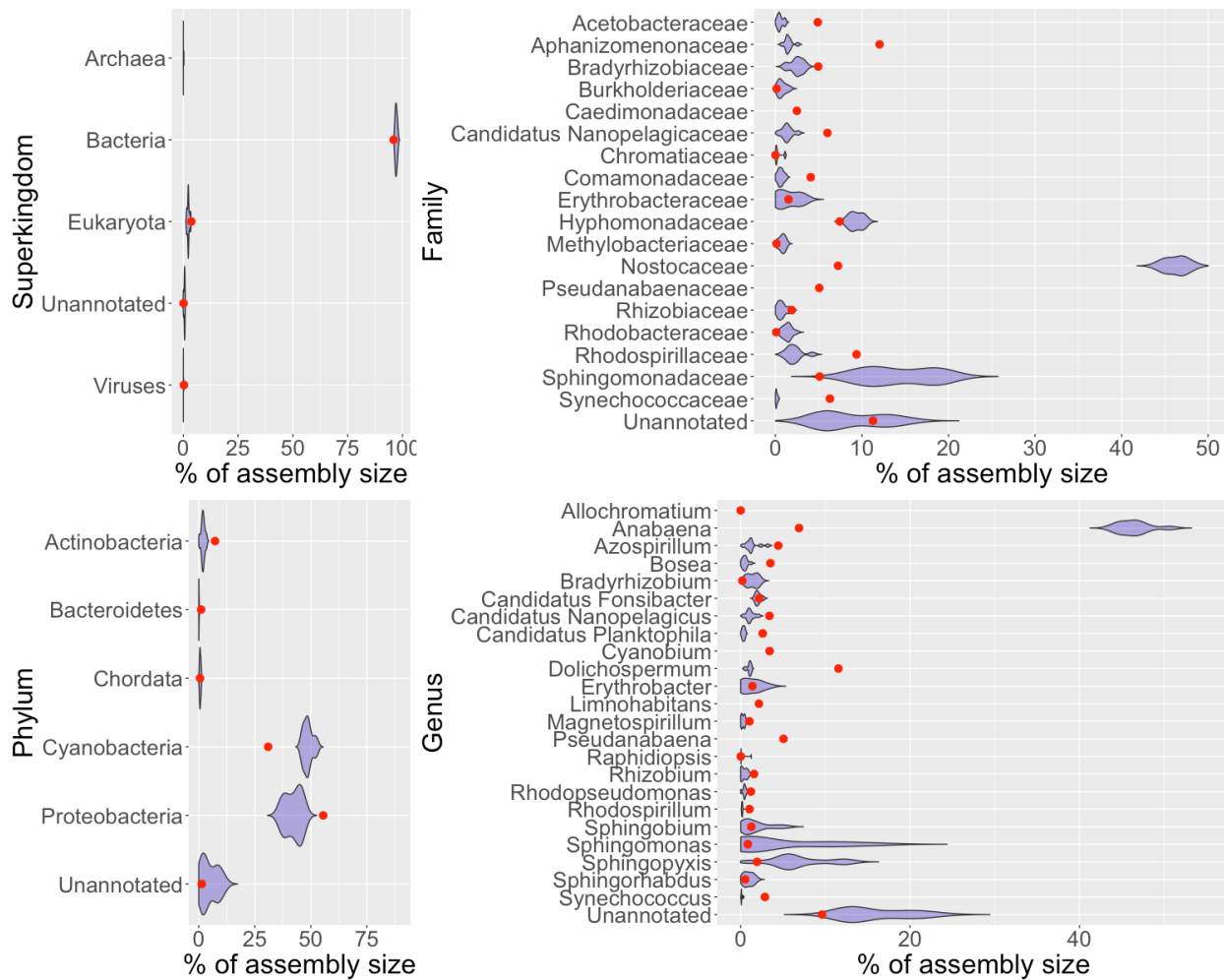
446



447

**Figure S5**. Comparison of assembly sizes for superkingdoms, phyla, families, and genera

between total assembly (dots) and Flongle-mimic subsets (violin plots) as annotated by BLAST.

Only taxonomic groups representing > 1% of assembly shown.

451

452

453

454

455

456

REFERENCES

457

458

459   Besemer, J., and M. Borodovsky, 1999 Heuristic approach to deriving models for gene finding.

460   Nucleic Acids Research 27(19): 3911-3920. https://doi.org/10.1093/nar/27.19.3911

461

462   Burton, A.S., S.E. Stahl, K.K John, M. Jain, S. Juul, et al., 2020 Off Earth identification of

463   bacterial populations using 16S rDNA nanopore sequencing. Genes 11: 76.

464

465   Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, et al., 2009 BLAST+:

466   architecture and applications. BMC Bioinformatics 10: 421.

467

468   Castro-Wallace, S.L., C.Y. Chiu, K.K. John, S.E. Stahl, K.H. Rubins *et al.,* 2017 Nanopore DNA

469   sequencing and genome assembly on the International Space Station. Scientific Reports 7:

470   18022. http://doi.org/10.1038/s41598-017-18364-0

471

472   Dietrich, A., 2006 Aesthetic issues for drinking water. Journal of Water and Health 4(S1): 11-16.

473

474   Dunlap, C.R., K.S. Sklenar, and L.J. Blake, 2015 A costly endeavor: addressing algae problems

475   in a water supply. Journal of the American Water Works Association 107(5): E255-E262.

476

477   Etheridge, S.M., 2010 Paralytic shellfish poisoning: Seafood safety and human health

478   perspectives. Toxicon 56(2): 108-122.

479

480   Gamez, T., 2018 The use of 23S ribotyping to detect harmful and nuisance phytoplankton in a

481   large, subtropical reservoir during an extended drought period (Unpublished thesis). Texas

482   State University, San Marcos, Texas.

483

484    Gilbert, P.M., D.M. Anderson, P. Gentien, E. Granéli, and K.G. Sellner, 2005 The global,

485    complex phenomena of harmful algal blooms. Oceanography 18(2): 136-147.

486

487    Gorham, T., E.D. Root, Y. Jia, C.K. Shum, and J. Lee, 2020. Relationship between

488    cyanobacterial bloom impacted drinking water sources and hepatocellular carcinoma incidence

489    rates. Harmful Algae 95: 101801.

490

491    Gowers, G.-O. F., O. Vince, J.-H. Charles, I. Klarenberg, T. Ellis, et al., 2019 Entirely off-grid

492    and solar-powered DNA sequencing of microbial communities during an ice cap traverse

493    expedition. Genes 10: 902. https://doi.org/10.3390/genes10110902

494

495    Griffith, A.W., and C.J. Gobler, 2020 Harmful algal blooms: a climate change co-stressor in

496    marine and freshwater ecosystems. Harmful Algae 91: 101590.

497

498    Hatfield, R.G., F.M. Batista, T.P. Bean, V.G. Fonseca, A. Santos, et al., 2020 The application of

499    nanopore sequencing technology to the study of dinoflagellates: a proof of concept study for

500    rapid sequence-based discrimination of potentially harmful algae. Frontiers in Microbiology 11:

501    844.

502

503    Hennon, G.M.M., and S.T. Dyhrman, 2020 Progress and promise of omics for predicting the

504    impacts of climate change on harmful algal blooms. Harmful Algae 91: 101587.

505

506    Huerta-Cepas, J., D. Szklarczyk, D. Heller, A. Hernández-Plaza, S.K. Forslund, et al., 2019

507    eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource

508  based on 5090 organisms and 2502 viruses. Nucleic Acids Research 47(D1): D309-D314.

509  https://doi.org/10.1093/nar/gky1085

510

511  Hoagland, P., and S. Scatasta, 2006 The economic effects of harmful algal blooms, in *Ecology*

512  *of Harmful Algae. Ecological Studies (Analysis and Synthesis) vol 189,* edited by E. Granéli and

513  J.T. Turner. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-540-32210-8_30

514

515  Johnson, S.S., E. Zaikova, D.S. Goerlitz, Y. Bai, and S.W. Tighe, 2017 Real-time DNA

516  sequencing in the Antarctic dry valleys using the Oxford Nanopore sequencer. Journal of

517  Biomolecular Techniques 28(1): 2-7. https://doi.org/10.7171/jbt.17-2801-009.

518

519  Katoh, K., and D.M. Standley, 2013 MAFFT multiple sequence alignment software version 7:

520  improvements in performance and usability. Molecular Phylogenetics and Evolution 30(4): 772-

521  780.

522

523  Kim, D., L. Song, F.P. Breitweiser, and S.L. Salzberg, 2016 Centrifuge: rapid and sensitive

524  classification of metagenomic sequences. Genome Research 26(12): 1721-1729.

525

526  Krehenwinkel, H., M. Wolf, J.Y. Lim, A.J. Rominger, W.B. Simison, et al., 2017 Estimating and

527  mitigating amplification bias in qualitative and quantitative arthropod metabarcoding. Scientific

528  Reports 7: 17668. https://doi.org/10.1038/s41598-017-17333-x

529

530  Kolmogorov, M., J. Yuan, Y. Lin, and P.A. Pevzner, 2019 Assembly of long, error-prone reads

531  using repeat graphs. Nature Biotechnology 37**:** 540-546. https://doi.org/10.1038/s41587-019-

532  0072-8

533

534     Legrand, B., J. Lesobre, J. Colombet, D. Latour, and M. Sabart, 2016 Molecular tools to detect

535     anatoxin-a genes in aquatic ecosystems: Toward a new nested PCR-based method. Harmful

536     Algae 58: 16-22.

537

538     Lepère, C., A. Wilmotte, and B. Meyer, 2000 Molecular diversity of *Microcystis* strains

539     (Cyanophyceae, Chroococcales) based on 16S rDNA sequences. Systematics and Geography

540     of Plants 70(2): 275-283

541

542     Li, H, 2018 Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34(18):

543     3094-3100. https://doi.org/10.1093/bioinformatics/bty191

544

545     Martí, J.M., and C.P. Garay, 2019 Not just BLAST nt: WGS database joins the party. *bioRxiv*

546     653592. https://doi.org/10.1101/653592 (Preprint posted June 4, 2019)

547

548     Menegon, M., C. Cantaloni, A. Rodriguez-Prieto, C. Centomo, A. Abdelfattah, et al., 2017 On

549     site DNA barcoding by nanopore sequencing. PLOS ONE *12*(10): e0184741.

550     https://doi.org/10.1371/journal.pone.0184741

551

552     Nguyen, L.-T., H.A. Schmidt, A. von Haeseler, and B.Q. Minh, 2015 IQ-TREE: A fast and

553     effective stochastic algorithm for estimating maximum likelihood phylogenies. Molecular Biology

554     and Evolution 32: 268-274. https://doi.org/10.1093/molbev/msu300

555

556     Pfeiffer, F., C. Gröber, M. Blank, K. Händler, M. Beyer, *et al.,* 2018 Systematic evaluation of

557     error rates and causes in short samples in next-generation sequencing. Scientific Reports 8:

558     10950. https://doi.org/10.1038/s41598-018-29325-6

559

560    Pomerantz, A., N. Peñafiel, A. Arteaga, L. Bustamante, F. Pichardo, et al., 2018 Real-time DNA

561    barcoding in a rainforest using nanopore sequencing: opportunities for rapid biodiversity

562    assessments and local capacity building. Gigascience 7: 1-14.

563

564    Quast, C., E. Pruesse, P. Yilmaz, J. Gerken, T. Schweer, et al., 2013 The SILVA ribosomal

565    RNA gene database project: improved data processing and web-based tools. Nucleic Acids

566    Research 41: D590-D596.

567
568
569    Rang, F.J., W.P. Kloosterman, and J. de Ridder, 2018 From squiggle to basepair: computational

570    approaches for improving nanopore sequencing read accuracy. Genome Biology 19: 90.

571    https://doi.org/10.1186/s13059-018-1462-9

572

573    Schafran, G.C, 2005 Reservoir Management Techniques to Improve Raw Water Quality.

574    Algal Metabolytes Workshop, December 9, 2005, Sarasota, FL.

575

576    Seppey M., M. Manni, and E.M. Zdobnov, 2019 BUSCO: Assessing Genome Assembly and

577    Annotation Completeness in *Gene Prediction. Methods in Molecular Biology vol 1962,* edited by

578    M. Kollmar. Humana, New York, NY. doi.org/10.1007/978-1-4939-9173-0_14 .

579

580    Sherwood, A.R., and G.G. Presting, 2007 Universal primers amplify a 23S rDNA plastid marker

581    in eukaryotic algae and cyanobacteria. Journal of Phycology 43(3): 605-608.

582    https://doi.org/10.1111/j.1529-8817.2007.00341.x

583

584    Stern, R.F., R.A. Andersen, I. Jameson, F.C. Küpper, M.A. Coffroth, et al., 2012 Evaluating the

585    ribosomal internal transcribed spacer (ITS) as a candidate dinoflagellate barcode marker. PLOS

586    ONE 7(8): e42780. https://doi.org/10.1371/journal.pone.0042780

587

588    Suurnäkki, S., G.V. Gomez-Saez, A. Rantala-Ylinen, J. Jokela, D.P. Fewer, et al., 2015

589    Identification of geosmin and 2-methylisoborneol in cyanobacteria and molecular detection

590    methods for the producers of these compounds. Water Research 68: 56-66.

591

592    Tsao, H.-W., A. Michinaka, H.-K. Yen, S. Giglio, P. Hobson, et al., 2014 Monitoring of geosmin

593    producing *Anabaena circinalis* using quantitative PCR. Water Research 49:416-425.

594    http://dx.doi.org/10.1016/j.watres.2013.10.028.

595

596    Vaser, R., I. Sović, N. Nagarajan, and M. Šikić, 2017 Fast and accurate de novo genome

597    assembly from long uncorrected reads. Genome Research *27*: 737-746.

598    https://doi.org/10.1101/gr.214270.116

599

600    Wacklin P., L. Hoffmann, and J. Komárek, 2009 Nomenclatural validation of the genetically

601    revised cyanobacterial genus *Dolichospermum* (Ralfs ex Bornet et Flahault) comb. nova. Fottea

602    9(1): 59-64.

603

604    Wick, R.R., L.M. Judd, and K.E. Holt, 2019 Performance of neural network basecalling tools for

605    Oxford Nanopore sequencing. Genome Biology 20: 129. https://doi.org/10.1186/s13059-019-

606    1727-y

607

608    Wick R.R., M.B. Schultz, J. Zobel, and K.E. Holt, 2015. Bandage: interactive visualisation of *de*

609    *novo* genome assemblies. Bioinformatics 31(20): 3350-3352.

610

611   Wongsurawat T., M. Nakagawa, O. Atiq, H.N. Coleman, P. Jenjaroenpun, et al., 2019 An

612   assessment of Oxford Nanopore sequencing for human gut metagenome profiling: A pilot study

613   of head and neck cancer patients. Journal of Microbiological Methods 166:

614   105739. https://doi.org/10.1016/j.mimet.2019.105739

615

616   Wood, D.E., J. Lu, and B. Langmead, 2019 Improved metagenomic analysis with Kraken2.

617   Genome Biology 20: 257.

618

619   Zhang, Y., X. Lin, T. Li, H. Li, L. Lin, et al., 2020 High throughput sequencing of 18S rRNA and

620   its gene to characterize a *Prorocentrum shikokuense* (Dinophyceae) bloom. Harmful Algae 94:

621   101809.

622

623   Zhu, W., A. Lomsadze, and M. Borodovsky, 2010 *Ab initio* gene identification in metagenomic

624   sequences. Nucleic Acids Research 38(12): e132.  https://doi.org/10.1093/nar/gkq275