

1 **Draft genome assemblies using sequencing reads from Oxford Nanopore Technology and**
2 **Illumina platforms for four species of North American killifish from the *Fundulus* genus**

3

4 Lisa K. Johnson [1,2], Ruta Sahasrabudhe [3], Tony Gill [1], Jennifer Roach [1], Lutz Froenicke
5 [3], C. Titus Brown [2], Andrew Whitehead* [1]

6

7 [1] Department of Environmental Toxicology, University of California, Davis

8 [2] Department of Population Health & Reproduction, School of Veterinary Medicine,

9 University of California, Davis

10 [3] DNA Technologies Core, Genome Center, University of California, Davis

11 *corresponding author: awhitehead@ucdavis.edu

12

13 **Abstract**

14

15 Draft *de novo* reference genome assemblies were obtained from four North American killifish
16 species (*Fundulus xenicus*, *Fundulus catenatus*, *Fundulus nottii*, and *Fundulus olivaceus*) using
17 sequence reads from Illumina and Oxford Nanopore Technologies' PromethION platforms. For
18 each species, the PromethION platform was used to generate 30-45x sequence coverage, and the
19 Illumina platform was used to generate 50-160x sequence coverage. Contig N50 values ranged
20 from 0.4 Mb to 2.7 Mb, and BUSCO scores were consistently above 90% complete using the
21 Eukaryota database. Draft assemblies and raw sequencing data are available for public use. We
22 encourage use and re-use of these data for assembly benchmarking and external analyses.

23

24 **Keywords:** long reads; Oxford Nanopore; killifish; genomes; genome assembly

25

26 **Background**

27

28 Sequencing and assembling large eukaryotic genomes is challenging [1–3]. Accuracy of
29 downstream analyses, such as variant calling and measuring gene expression, depends heavily on
30 a high-quality reference genome [4]. Fortunately, the cost of generating whole genome sequence
31 data is dropping, making it easier for individual labs rather than large consortiums to generate
32 assemblies for organisms without reference genomes [3,5,6]. Single-molecule long read nucleic
33 acid sequencing technology from Oxford Nanopore Technologies (ONT), which has been
34 commercially available since 2014 [7], has been shown to improve the contiguity of reference
35 assemblies [8] and reveal “dark regions” that were previously camouflaging genes [9]. The
36 lengths of the sequencing reads generated using this technology are limited only by the size of
37 the fragments in the extracted DNA sample [10]. The promise of more complete reference
38 assemblies is especially important for the accuracy of comparative evolutionary genomics
39 studies, as assembly fragments lead to errors in downstream synteny analyses [11], as well as

40 SNP calling and identification of transcript features (splice junctions and exons) for
41 quantification.

42

43 Despite high error rates ~5% [12] relative to Illumina short reads ~0.25% [13] and the relatively
44 recent availability of ONT data, there have been a flurry of studies using this sequencing
45 technology. Small genomes from bacteria and viruses appear to be ideal for sequencing on the
46 ONT MinION platform [12]. The portable nature of the technology makes it appealing as a
47 resource for teaching [14,15], working in remote locations [16–18] and for investigating viral
48 outbreak public health emergencies [19–21]. However, despite the demonstrated ability to
49 achieve yields >6.5 Gb per flow cell [22], the MinION platform can be prohibitively expensive
50 for sequencing larger eukaryotic genomes. For example, 39 flow cells yielded 91.2 Gb of
51 sequence data (~30x coverage) of the human genome [23]. Sequencing of the wild tomato
52 species, *Solanum pennellii* across thirty-one flow cells yielded 110.96 Gb (~100x coverage) with
53 some flow cells yielding >5Gb [24]. By contrast, following the 2018 beta release of the ONT
54 PromethION platform, which has a higher density of nanopore channels, five flow cells were
55 used to yield >250 Gb (~80x coverage) of the human genome [25]. PromethION data combined
56 with Hi-C long-range mapping data from human samples produced a genome genome assembly
57 with a scaffold N50 of 56.4 Mbp [26].

58

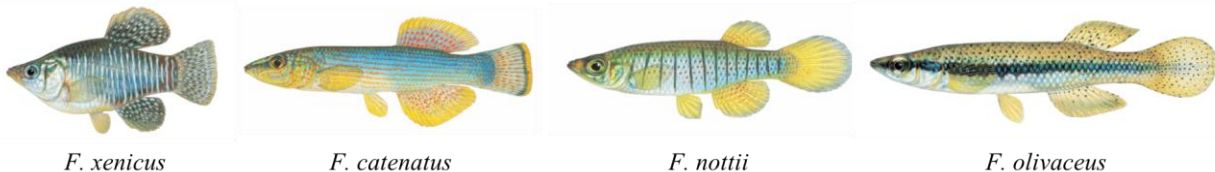
59 The combination of long read sequencing data from ONT MinION and short read sequencing
60 data from Illumina has been used to improve the quality of reference genomes [27–30]. In one
61 approach, short read assembly scaffolds have been improved with the addition of long reads. The
62 Murray cod genome (640-669 Mb in size) was improved by combining low coverage 804 Mb of
63 long reads ONT data from just one MinION flow cell with 70.6 Gb of Illumina data from both
64 HiSeq and MiSeq; the assembly scaffold N50 increased from 33,442 bp (Illumina only) to
65 52,687 bp with ONT and Illumina combined [31]. The clownfish genome (791 to 794 Mb in
66 size) was improved by including 8.95 Gb of ONT MinION reads; the scaffold N50 increased
67 from 21,802 bp (Illumina only) to 401,715 bp with ONT and Illumina combined [27]. Recently,
68 a new approach is available with racon [32] and/or pilon [33] consensus building tools, which
69 uses Illumina data to “polish” contigs from ONT-only assemblies. Polishing corrects single
70 nucleotide base differences, fills gaps, and identifies local mis-assemblies [33]. This approach
71 has been shown to improve the BUSCO score from <1% with the ONT assembly alone to >95%
72 complete after polishing with Illumina reads [28].

73

74 In this study, we explored whether the ONT PromethION sequencing technology could be
75 appropriate for generating initial draft reference genomes for four species of North American
76 killifish belonging to the *Fundulus* genus. *Fundulus* is a comparative evolutionary model system
77 for studying repeated genomic divergence between marine and freshwater species. *Fundulus*
78 killifish have a cosmopolitan geographic distribution across North America. These small
79 cyprinodontiform fish have evolved to occupy a wide range of osmotic niches, including marine,

80 estuarine, and freshwater [34]. Estuarine and coastal *Fundulus* are euryhaline, insofar as they can
81 adjust their physiologies to tolerate a very wide range of salinities. In contrast, freshwater species
82 are stenohaline: they tolerate a much narrower range of salinities [34,35]. Freshwater clades are
83 derived from marine clades, and radiation into freshwater has occurred multiple times
84 independently within the genus. This makes *Fundulus* unusual, because most large clades of
85 fishes are either exclusively marine or exclusively freshwater. Therefore, species of closely-
86 related killifish in the *Fundulus* genus serve as a unique comparative model system for
87 understanding the genomic mechanisms that contribute to evolutionary divergence and
88 convergence of osmoregulatory processes, which is important for understanding how species will
89 evolve to face fluctuating salinities in future climate change scenarios [36]. The Atlantic killifish,
90 *Fundulus heteroclitus* has been a well-described physiological model organism for investigating
91 the functional basis of, and evolution of, physiological resilience to temperature, salinity,
92 hypoxia, and environmental pollution [34,37–39], with an available genome from the Atlantic
93 killifish, *Fundulus heteroclitus* [40]. However, we do not currently have any genomes from other
94 *Fundulus* killifish, particularly from those occupying freshwater habitats.

95
96 Here, we report the collection of whole genome sequencing data using both ONT PromethION
97 and Illumina platforms from four killifish species without previously-existing sequencing data
98 (Figure 1): *Fundulus xenicus* (formerly *Adinia xenica*) [41], *Fundulus catenatus*, *Fundulus nottii*,
99 and *Fundulus olivaceus*. *F. xenicus* is euryhaline and occupies coastal and estuarine habitats,
100 while the other species (*F. catenatus*, *F. nottii*, *F. olivaceus*) are stenohaline and occupy
101 freshwater habitats.



103 Figure 1. Four *Fundulus* killifish (left to right): the marine diamond killifish *Fundulus xenicus*;
104 the northern studfish, *Fundulus catenatus* (south central United States); the freshwater bayou
105 topminnow, *Fundulus nottii*; and the freshwater blackspotted topminnow, *Fundulus olivaceus*.
106 (drawings used with permission from the artist, Joseph R. Tomelleri).

107 108 **Methods and Results**

109
110 Live field-caught individuals of each fish species were shipped to UC Davis and kept at their
111 native salinities in an animal holding facility, maintained according to University of California
112 IACUC standards. *F. catenatus* and *F. olivaceus* were collected from the Gasconade River, MO
113 (latitude/longitude coordinates 37.879/-91.795 and 37.19/-92.56, respectively), *F. nottii* was
114 collected from Walls Creek, MS (31.154433/-89.245381), and *F. xenicus* was collected from
115 Graveline Bayou, MS (30.368756/-88.719329). High molecular weight (hmw) DNA was
116 extracted from fresh tissue for *F. nottii* and *F. xenicus*, and from frozen tissue for *F. catenatus*

117 and *F. olivaceus*. For *F. catenatus* and *F. olivaceus*, tissues were dissected and frozen in liquid
118 nitrogen and stored at -80 °C until samples were prepared for hmw DNA extraction. With the
119 exception of *F. olivaceus*, each assembly consisted of sequencing one tissue sample from one
120 individual. For *F. olivaceus*, Illumina data were collected from DNA extracted from one
121 individual while the ONT PromethION data were collected from another individual (frozen
122 tissue).

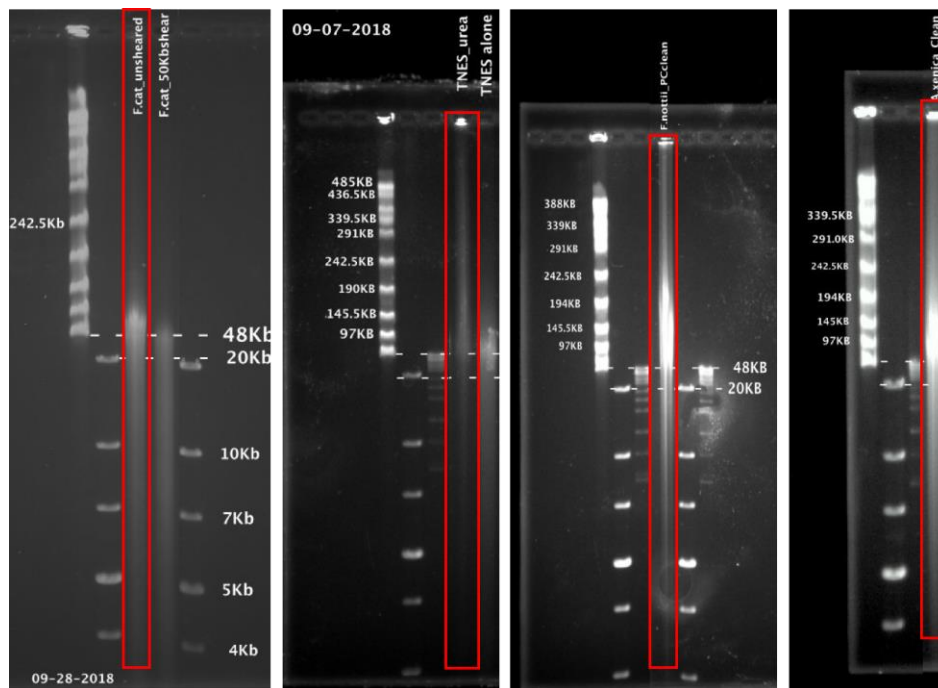
123

124 *DNA extractions*

125

126 Whole fish heads were used for hmw DNA extractions. Agilent's Genomic DNA Isolation kit
127 (Catalog #200600) was used to extract DNA from fresh tissues from *F. xenicus* and *F. nottii*. For
128 *F. catenatus* and *F. olivaceus*, both the ultra-long read sequencing protocol from [42] (which
129 included Tissue lysis buffer with Tris, NaCl, EDTA, SDS and Proteinase K followed by
130 phenol:chloroform extraction), as well as the Qiagen "DNA purification from tissue using the
131 Genra puregene Tissue Kit" (p. 39) were used, and were found to be similar to the Agilent kit.
132 Precipitated DNA was difficult to re-dissolve; therefore, additional phenol:chloroform cleanup
133 steps were required after extractions. We found that adding urea to the lysis buffer helped the
134 precipitated DNA pellet to be less fragile and go into solution easier [43]. Prior to library
135 preparation, hmw DNA from *F. nottii* and *F. olivaceus* (PromethION) was sheared to 50 kb in an
136 effort to improve the ligation enzyme efficiency, resulting in fragments in the 50-70 kb range.
137 Field inversion gels were used to visualize hmw DNA (Figure 2).

138



139

140 Figure 2. Field inversion gels with red boxes showing samples sequenced (in order from left to
141 right: *F. catenatus* (sheared vs. unsheared), *F. olivaceus*, *F. nottii*, *F. xenicus*). DNA was

142 extracted from fresh tissues for *F. xenicus* and *F. nottii*, and from frozen tissue for *F. catenatus*
143 and *F. olivaceus*.

144

145 *ONT sequencing*

146

147 Libraries for ONT PromethION sequencing were prepared using the ligation sequencing kit
148 (SQK-LSK109) following the manufacturer's instructions. ONT PromethION sequencing data
149 were collected from all four species on an alpha-beta instrument through the early release
150 program at the University of California, Davis DNA Technologies Core facility (Davis, CA
151 USA). One species was sequenced per flow cell. Base-calling was done onboard the
152 PromethION instrument (Oxford Nanopore Technologies, UK). For the *F. xenicus* run, lambda
153 phage (DNA CS) was spiked-in as a positive control.

154

155 *Illumina Sequencing*

156

157 With the exception of *F. olivaceus*, each individual hmw DNA sample used for the ONT library
158 was also used for Illumina library preparation using the Nextera Index Kit (FC-121-1012). For
159 three species, Illumina data were multiplexed across two PE150 lanes on an Illumina HiSeq 4000
160 and demultiplexed by Novogene (Sacramento, CA USA). For *F. olivaceus*, PE150 Illumina
161 NovaSeq reads from one flow cell (2 lanes) were graciously provided by the Texas A&M
162 Agrilife Research Sequencing Facility (College Station, TX USA).

163

164 **Data Description**

165

166 Whole genome sequencing data from individuals of four killifish species collected from ONT
167 PromethION (Table 1) and Illumina (NovaSeq and HiSeq 4000) (Table 2) were deposited in the
168 European Nucleotide Archive (ENA) under the study accession PRJEB29136. Deposited raw
169 data are untrimmed and unfiltered. Reads corresponding to lambda phage were filtered from
170 ONT PromethION data using the NanoLyse program from NanoPack (version 1.1.0; [44]).
171 Porechop (version 0.2.3) was used to remove residual ONT adapters and NanoFilt (version 2.2.0;
172 [44]) was used to filter reads with an average quality score >Q5. After filtering and adapter
173 trimming, ONT data from the PromethION ranged from 30-45x coverage for each species.
174 NanoPlot (version 1.10.0; [44]) was used for visualization of ONT read qualities.

175

176

177

178

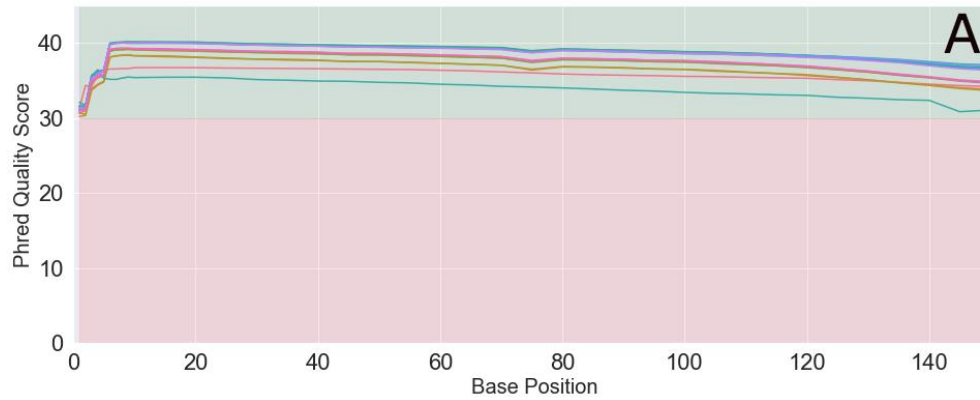
179

Species	Bases called (Gb)	Cov. (x)	Average read length	Reads N50	Q>5 Bases called (Gb)	Q>5 Avg. read length	ONT signal Accession	ONT fastq Accession
<i>F. xenicus</i>	38.5	35	2,449	5,733; n = 1,373,426	36.42	2,699	ERR3385273	ERR3385269
<i>F. nottii</i>	33.4	30.4	6,480	12,995; n = 700,534	31.06	7,548	ERR3385275	ERR3385271
<i>F. catenatus</i>	40.3	36.6	1,699	3,439; n = 2,687,295	34.28	2,021	ERR3385274	ERR3385270
<i>F. olivaceus</i>	50.1	45.5	4,595	11,670; n = 987,921	45.97	5,365	ERR3385276	ERR3385272

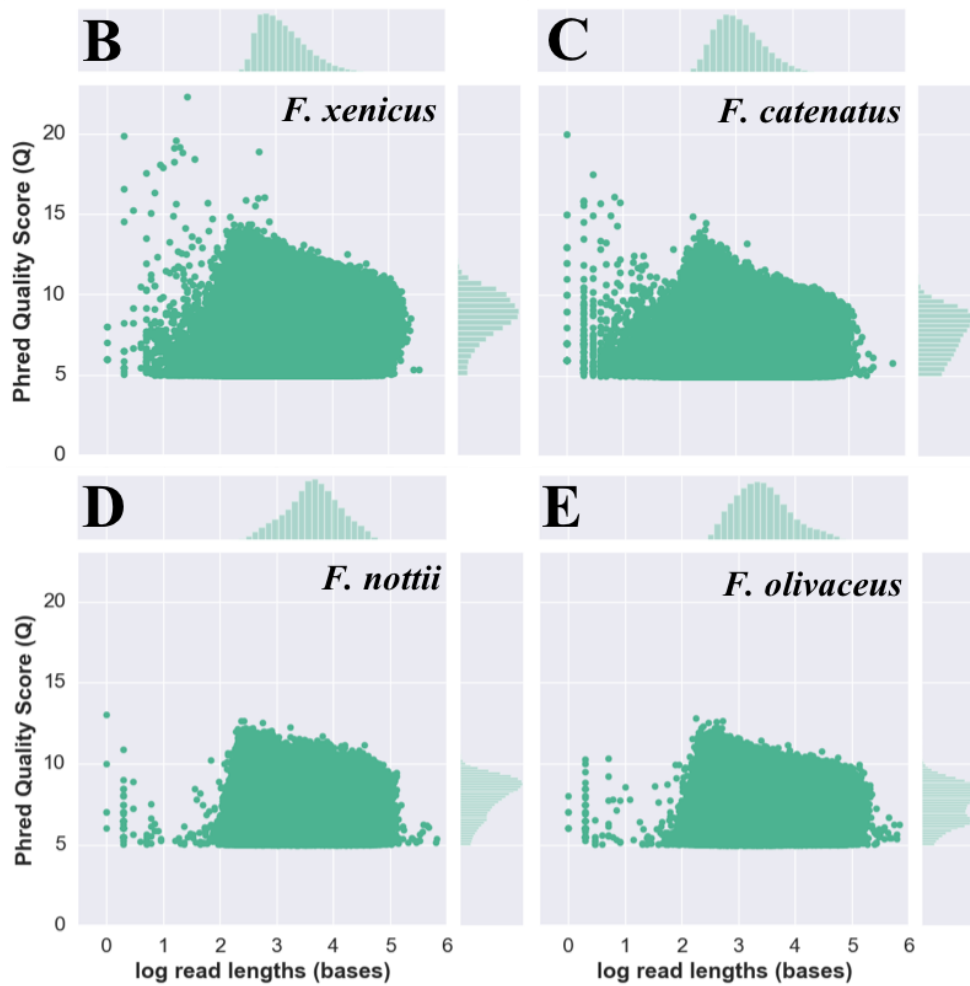
180 Table 1. ONT data collected from each species. Coverage assumes the genome size of each
 181 species is 1.1 Gb, measured for *F. heteroclitus* (Reid et al. 2017). Untrimmed reads were
 182 deposited in the ENA under study PRJEB29136. Q>5 reads were used in subsequent assemblies.
 183

184 Average quality scores for all Illumina data were consistently above Q30 (Figure 3A). Residual
 185 Nextera adapters and bases with low quality scores were removed from Illumina reads using
 186 Trimmomatic PE (version 0.38) with conservative parameters, which included removing bases
 187 from each read with a quality score below Q2 and required a minimum read length of 25 bases
 188 each [45].
 189

190 For *F. xenicus* and *F. catenatus*, ONT read qualities ranged from Q5 (minimum cutoff) to Q14
 191 with read lengths generally ranging from 10 bp to 100kb (Figure 3B,C). For *F. nottii* and *F.*
 192 *olivaceus*, ONT read qualities ranged from Q5 (minimum cutoff) to Q13 with read lengths
 193 ranging from 100 bp to 100kb (Figure 3D,E).



194



195

196 Figure 3. A) Quality score profiles for PE Illumina reads from *F. xenicus*, *F. catenatus*, *F. nottii*
197 and *F. olivaceus*. For Illumina data, phred quality scores were consistently above Q30 across all
198 reads. Average read quality scores (Q score) vs. read lengths for ONT PromethION from B) *F.*
199 *xenicus*, C) *F. catenatus*, D) *F. nottii*, E) *F. olivaceus*.

200

201

Species	Platform	Reads (M)	Coverage (x)	FASTQ Accessions
<i>F. xenicus</i>	Illumina HiSeq	327.5	89.3	ERR3385278 ERR3385279
<i>F. nottii</i>	Illumina HiSeq	197	53.7	ERR3385282 ERR3385283
<i>F. catenatus</i>	Illumina HiSeq	316.5	86.3	ERR3385280 ERR3385281
<i>F. olivaceus</i>	Illumina NovaSeq	601.9	164	ERR3385284 ERR3385285

202 Table 2. Illumina data collected were all paired-end PE 150 reads. Coverage assumes 1.1 Gb
203 genome size measured for *F. heteroclitus* [40].

204

205 Draft Assemblies

206

207 As a comparison with assemblies using long ONT read data, Illumina data alone were assembled
208 using ABySS version 2.1.5. While the BUSCO scores were consistently above 50%
209 completeness [46], the number of contigs and contig N50 lengths of the Illumina-only assemblies
210 were not acceptable for downstream use (Table 3).

211

Species	Bases in the Illumina-only assembly	N contigs	Avg length	Largest contig	N50	Illumina-only BUSCO C CS/CD/F/M
<i>F. xenicus</i>	1,283,257,056	5,195,861	246.98	71,596	2,571; n = 107,350	57.1% 56.4/0.7/33.3/9.6
<i>F. catenatus</i>	1,205,429,912	3,989,534	302.15	70,870	3,629; n = 80,839	53.8% 52.8/1.0/36.0/10.2
<i>F. nottii</i>	1,167,835,004	3,875,693	301.32	92,540	3,740; n = 72810	62.7% 61.7/1.0/27.4/9.9
<i>F. olivaceus</i>	1,252,948,998	4,509,089	277.87	70,765	3,670; n = 77136	65.7% 64.0/1.7/25.1/9.2

212 Table 3. Statistics for Illumina-only assemblies using ABySS (version 2.1.5) for each species.

213 The BUSCO Eukaryota database (303 genes) was used to evaluate the completeness of each

214 assembly [46]. BUSCO numbers reported are percentage complete (C) followed by the
 215 percentages of complete single-copy (CS), complete duplicated (CD), fragmented (F), missing
 216 (M) out of 303 genes.

217
 218 The ONT-only assemblies using the fuzzy de Bruijn graph assembler, wtdbg2 (version 2.3; [47])
 219 had high contig N50 but low complete matches with the BUSCO Eukaryota database (Table 4).
 220 The assembler wtdbg2 took an average of 6.1 hours per assembly and required 59 GB RAM. The
 221 polishing tool pilon required an average of 65.99 hours and used 1.61 TB RAM. Following
 222 polishing with Illumina data using the pilon software tool version 1.23 [33], the BUSCO
 223 Eukaryota completeness scores increased to consistently greater than 90% (Table 4). A full table
 224 of BUSCO metrics can be found in Supplemental Table 1. Assemblies were deposited in the
 225 Open Science Framework (OSF) repository, <https://doi.org/10.17605/osf.io/zjv86> and in zenodo
 226 record, <https://doi.org/10.5281/zenodo.3251033>.

227

Species	Contigs	Contig N50	Assembly size (bases)	Complete BUSCO after wtdbg2 ONT-only	Complete BUSCO after pilon polishing with Illumina
				C CS/CD/F/M	C CS/CD/F/M
<i>F. xenicus</i>	5,621	888,041; n = 325	1,075,031,690	10.2% 10.2/0.0/11.6/78.2	90.5% 87.5/3.0/3.0/6.5
<i>F. nottii</i>	2,242	2,701,963; n = 95	1,081,276,623	28.4% 11.2/0.0/22.1/66.7	94.4% 92.1/2.3/1.0/4.6
<i>F. catenatus</i>	5,854	436,102; n = 780	1,163,592,740	11.2% 28.4/0.0/24.4/47.2	90.4% 88.4/2.0/2.6/7.0
<i>F. olivaceus</i>	2,622	2,669,230; n = 105	1,198,526,423	23.4% 23.4/0.0/25.7/50.9	92.1% 89.8/2.3/1.3/6.6

228 Table 4. ONT PromethION assemblies using the wtdbg2 version 2.3 assembler [47] followed by
 229 polishing with pilon version 1.23 [33]. Of interest is the dramatic improvement of the complete
 230 BUSCO metric after polishing with pilon. BUSCO numbers reported are percentage complete
 231 (C) followed by the percentages of complete single-copy (CS), complete duplicated (CD),
 232 fragmented (F), missing (M) out of the 303 genes in the BUSCO Eukaryota database [46].

233

234

235

236 Discussion

237

238 In this study, we collected 30-45x coverage of ONT data in combination with 50-160x coverage
239 of Illumina PE150 sequencing data and generated draft genome assemblies for four species of
240 *Fundulus* killifish. For these assemblies, the combination of ONT and Illumina data allowed us
241 to generate highly contiguous assemblies with acceptable BUSCO results. The assemblies
242 generated by ONT data alone were not acceptable for use because of the low BUSCO results,
243 due to the high rate of ONT sequence errors. Polishing the ONT assemblies with the Illumina
244 data did not improve contiguity of the assemblies, but served to correct errors, shown by the
245 large boost in BUSCO scores relative to the ONT assemblies alone.

246

247 The qualities of the ONT data appeared to make a difference in the contig N50 metrics of the
248 assemblies. *F. nottii* and *F. olivaceus* both had contig N50 >2Mb, while assemblies from *F.*
249 *xenicus* and *F. catenatus* had contig N50 <1Mb. *F. xenicus* and *F. catenatus* had shorter average
250 read lengths and reads N50, on average, compared to *F. nottii* and *F. olivaceus*. *F. nottii*, which
251 had the lowest yield, had higher average read lengths and reads N50 compared to the other
252 species. *F. olivaceus*, which had the highest yield, also had a high reads N50 and average read
253 length. Therefore, when generating ONT data for draft genome assemblies, it might matter more
254 to have a lower yield of longer reads than a higher yield of shorter reads.

255

256 While the ONT data collected were sufficient for genome assembly of these organisms, it is
257 worth noting that our yields were lower than those advertised on the ONT website (~100 Gb
258 from a single PromethION flow cell) (<https://nanoporetech.com/products/promethion>, accessed
259 06/12/2019), and read length N50 was shorter than expected based on DNA gel analysis. We
260 observed that our samples were consistently not using pores as efficiently on the PromethION
261 compared with other runs with samples isolated from human or mammalian samples. The
262 reasons for this are not fully understood, but this could be due to brittle property of our hmw
263 DNA.

264

265 The Vertebrate Genome Project (VGP 2018) lists standards for *de novo* long-range genome
266 assembly that include PacBio long reads, 10x linked Illumina reads, Hi-C chromatin mapping
267 and Bionano Genomics optical maps. These four types of data each have a high cost of
268 generation as well as associated analysis time and computational costs. While chromatin capture
269 and Hi-C methods produce chromosome-level assemblies of very high quality [48–51], this can
270 significantly increase the cost of the genome sequencing project. Here, we report the pairing of
271 high-quality short Illumina reads with error-prone long reads generated from the ONT
272 PromethION platform to generate a draft assembly at a minimum cost. The qualities of these
273 assemblies are not as high as compared to the standards recommended by VGP (2018) with the
274 3.4.2.QV40 phased metric, which requires the assembly to be haplotype phased with a minimum
275 contig N50 of 1 million bp (1Mb), scaffold N50 of 10Mb, 90% of the genome assembled into

276 chromosomes and a base quality error of Q40, (VGP 2019). However, for many research
277 purposes these assemblies are sufficient, considering the low cost and that we have a high-
278 quality reference genome assembly for another species within the genus [40]. For *F. olivaceus*
279 and *F. nottii*, draft assemblies using wtdbg2 [47] and pilon polishing with Illumina data [33] had
280 contig N50 >1 Mb, which meets the minimum requirements for assemblies in downstream
281 synteny analyses [11].

282

283 New software tools and methods for base-calling, assembling and analyzing noisy ONT long
284 reads are being developed at a fast rate [52,53]. Because of this fast pace of software tool
285 development for ONT data, standard operating procedures are not available. While we use these
286 assemblies for their intended purpose of downstream comparative evolutionary analyses, the raw
287 data are shared here with the intent that others may use them for tool development and as new
288 workflow pipelines, algorithms, tools, and best practices emerge.

289

290 **Conclusions**

291

292 Sequencing data from the ONT PromethION and Illumina platforms combined can contribute to
293 assemblies of eukaryotic vertebrate genomes (>1 Gb). These sequencing data from wild-caught
294 individuals of *Fundulus* killifish species are available for use with tool development and
295 workflow pipelines. Ongoing work from our group is comparing genomic content between
296 *Fundulus* species to address questions about evolutionary mechanisms of osmoregulatory
297 divergence.

298

299 **Data re-use potential**

300

301 We encourage use and re-use of these data for external analyses. This collection of whole
302 genome sequencing data from the PromethION and Illumina platforms originates from wild-
303 caught individuals of closely-related *Fundulus* killifish species, obtained for the purpose of
304 downstream evolutionary genomic comparative analyses. These data add to the growing set of
305 public data available from ONT PromethION sequencing platform [25,54] which can be used for
306 developing base-calling and assembly algorithms with this type of data.

307

308 **Availability of supporting data and materials**

309

310 Raw data are available in the ENA under study PRJEB29136. Draft assembly data products and
311 quality assessment reports are available in the OSF repository:

312 <https://doi.org/10.17605/osf.io/zjv86> and zenodo: <https://doi.org/10.5281/zenodo.3251033>.

313 Scripts used for this analysis workflow are available at:

314 https://github.com/johnsolk/ONT_Illumina_genome_assembly

315

316

317 **List of abbreviations**

318

319 BUSCO = Benchmarking Universal Single-Copy Orthologs

320 ENA = European Nucleotide Archive

321 hmw DNA = high molecular weight DNA

322 ONT = Oxford Nanopore Technologies

323 OSF = Open Science Framework

324 PE = paired end

325 VGP = Vertebrate Genome Project

326

327 **Declarations**

328

329 *Ethical Approval*

330 UC Davis IACUC protocol #17221

331

332 *Consent for publication*

333 Not applicable.

334

335 *Competing Interests*

336 The authors declare that they have no competing interests.

337

338 *Funding*

339 Gordon and Betty Moore Foundation to CTB under award number GBMF4551. IU-TACC

340 Jetstream and PSC Bridges XSEDE allocations TG-BIO160028 and TG-MCB190015 to LKJ.

341

342 **Author's Contributions**

343

344 Sample extractions and library preparations were done by LKJ, RS, TG, JR. Project advising by

345 CTB and AW. Manuscript writing and editing by LKJ, RS, TG, JR, LF, CTB, and AW.

346

347 **Acknowledgements**

348

349 We thank Dr. David Duvernell at Missouri University of Science & Technology and Dr. Jacob

350 Schaefer at the University of Southern Mississippi for generously collecting and sending fish. A

351 special thank you goes to Dr. Charlie Johnson and Dr. Richard Metz at Texas A&M University

352 Agrilife Research Sequencing Facility for contributing Illumina NovaSeq data from *Fundulus*

353 *olivaceus*. Thanks to the instructors and participants at PoreCamp USA (June 2017) for their

354 helpful advice.

355 **References**

- 356 1. Mardis E, McPherson J, Martienssen R, Wilson RK, McCombie WR. What is finished, and
357 why does it matter. *Genome Res* [Internet]. 2002;12:669–71. Available from:
358 <http://dx.doi.org/10.1101/gr.032102>
- 359 2. Baker M. De novo genome assembly: what every biologist should know. *Nat Methods*
360 [Internet]. Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights
361 Reserved.; 2012;9:333. Available from: <https://doi.org/10.1038/nmeth.1935>
- 362 3. Eklom R, Wolf JBW. A field guide to whole-genome sequencing, assembly and annotation.
363 *Evol Appl* [Internet]. 2014;7:1026–42. Available from: <http://dx.doi.org/10.1111/eva.12178>
- 364 4. Stemple DL. So, you want to sequence a genome. *Genome Biol* [Internet]. 2013;14:128.
365 Available from: <http://dx.doi.org/10.1186/gb-2013-14-7-128>
- 366 5. Li F-W, Harkess A. A guide to sequence your favorite plant genomes. *Appl Plant Sci*
367 [Internet]. 2018;6:e1030. Available from: <http://dx.doi.org/10.1002/aps3.1030>
- 368 6. Dominguez Del Angel V, Hjerde E, Sterck L, Capella-Gutierrez S, Notredame C, Vinnere
369 Pettersson O, et al. Ten steps to get started in Genome Assembly and Annotation. *F1000Res*
370 [Internet]. 2018;7. Available from: <http://dx.doi.org/10.12688/f1000research.13598.1>
- 371 7. Ip CLC, Loose M, Tyson JR, de Cesare M, Brown BL, Jain M, et al. MinION Analysis and
372 Reference Consortium: Phase 1 data release and analysis. *F1000Res* [Internet]. 2015;4:1075.
373 Available from: <http://dx.doi.org/10.12688/f1000research.7201.1>
- 374 8. Tyson JR, O’Neil NJ, Jain M, Olsen HE, Hieter P, Snutch TP. MinION-based long-read
375 sequencing and assembly extends the *Caenorhabditis elegans* reference genome. *Genome Res*
376 [Internet]. 2018;28:266–74. Available from: <http://dx.doi.org/10.1101/gr.221184.117>
- 377 9. Ebbert MTW, Jensen TD, Jansen-West K, Sens JP, Reddy JS, Ridge PG, et al. Systematic
378 analysis of dark and camouflaged genes reveals disease-relevant genes hiding in plain sight.
379 *Genome Biol* [Internet]. 2019;20:97. Available from: <http://dx.doi.org/10.1186/s13059-019-1707-2>
- 381 10. Laver T, Harrison J, O’Neill PA, Moore K, Farbos A, Paszkiewicz K, et al. Assessing the
382 performance of the Oxford Nanopore Technologies MinION. *Biomol Detect Quantif* [Internet].
383 2015;3:1–8. Available from: <http://dx.doi.org/10.1016/j.bdq.2015.02.001>
- 384 11. Liu D, Hunt M, Tsai IJ. Inferring synteny between genome assemblies: a systematic
385 evaluation. *BMC Bioinformatics* [Internet]. 2018;19:26. Available from:
386 <http://dx.doi.org/10.1186/s12859-018-2026-4>
- 387 12. Tyler AD, Mataseje L, Urfano CJ, Schmidt L, Antonation KS, Mulvey MR, et al. Evaluation
388 of Oxford Nanopore’s MinION Sequencing Device for Microbial Whole Genome Sequencing
389 Applications. *Sci Rep* [Internet]. 2018;8:10931. Available from:
390 <http://dx.doi.org/10.1038/s41598-018-29334-5>

- 391 13. Pfeiffer F, Gröber C, Blank M, Händler K, Beyer M, Schultze JL, et al. Systematic
392 evaluation of error rates and causes in short samples in next-generation sequencing. *Sci Rep*
393 [Internet]. 2018;8:10950. Available from: <http://dx.doi.org/10.1038/s41598-018-29325-6>
- 394 14. Zeng Y, Martin CH. Oxford Nanopore sequencing in a research-based undergraduate course
395 [Internet]. *bioRxiv*. 2017 [cited 2019 Jun 20]. p. 227439. Available from:
396 <https://www.biorxiv.org/content/10.1101/227439v1>
- 397 15. Zaaier S, Columbia University Ubiquitous Genomics 2015 class, Erlich Y. Using mobile
398 sequencers in an academic classroom. *Elife* [Internet]. 2016;5. Available from:
399 <http://dx.doi.org/10.7554/eLife.14258>
- 400 16. Ducluzeau A-L, Tyson JR, Collins RE, Snutch TP, Hassett BT. Genome Sequencing of Sub-
401 Arctic Mesomycetozoean *Sphaeroforma sirkka* Strain B5, Performed with the Oxford Nanopore
402 minION and Illumina HiSeq Systems. *Microbiol Resour Announc* [Internet]. 2018;7. Available
403 from: <http://dx.doi.org/10.1128/MRA.00848-18>
- 404 17. Pomerantz A, Peñafiel N, Arteaga A, Bustamante L, Pichardo F, Coloma LA, et al. Real-time
405 DNA barcoding in a rainforest using nanopore sequencing: opportunities for rapid biodiversity
406 assessments and local capacity building. *Gigascience* [Internet]. 2018;7. Available from:
407 <http://dx.doi.org/10.1093/gigascience/giy033>
- 408 18. Boykin LM, Ghalab A, De Marchi BR, Savill A, Wainaina JM, Kinene T, et al. Real time
409 portable genome sequencing for global food security [Internet]. *bioRxiv*. 2018 [cited 2019 Jun
410 20]. p. 314526. Available from: <https://www.biorxiv.org/content/10.1101/314526v2>
- 411 19. Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, et al. Real-time, portable
412 genome sequencing for Ebola surveillance. *Nature* [Internet]. 2016;530:228–32. Available from:
413 <http://dx.doi.org/10.1038/nature16996>
- 414 20. Quick J, Grubaugh ND, Pullan ST, Claro IM, Smith AD, Gangavarapu K, et al. Multiplex
415 PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly
416 from clinical samples. *Nat Protoc* [Internet]. 2017;12:1261–76. Available from:
417 <http://dx.doi.org/10.1038/nprot.2017.066>
- 418 21. Kafetzopoulou LE, Pullan ST, Lemey P, Suchard MA, Ehichioya DU, Pahlmann M, et al.
419 Metagenomic sequencing at the epicenter of the Nigeria 2018 Lassa fever outbreak. *Science*
420 [Internet]. 2019;363:74–7. Available from: <http://dx.doi.org/10.1126/science.aau9343>
- 421 22. Schalamun M, Nagar R, Kainer D, Beavan E, Eccles D, Rathjen JP, et al. Harnessing the
422 MinION: An example of how to establish long-read sequencing in a laboratory using challenging
423 plant tissue from *Eucalyptus pauciflora*. *Mol Ecol Resour* [Internet]. 2019;19:77–89. Available
424 from: <http://dx.doi.org/10.1111/1755-0998.12938>
- 425 23. Jain M, Koren S, Quick J, Rand AC, Sasani TA, Tyson JR, et al. Nanopore sequencing and
426 assembly of a human genome with ultra-long reads [Internet]. *bioRxiv*. 2017 [cited 2019 Jun 20].
427 p. 128835. Available from: <https://www.biorxiv.org/content/10.1101/128835v1>

- 428 24. Schmidt MH-W, Vogel A, Denton AK, Istace B, Wormit A, van de Geest H, et al. De Novo
429 Assembly of a New *Solanum pennellii* Accession Using Nanopore Sequencing. *Plant Cell*
430 [Internet]. 2017;29:2336–48. Available from: <http://dx.doi.org/10.1105/tpc.17.00521>
- 431 25. De Coster W, De Rijk P, De Roeck A, De Pooter T, D’Hert S, Strazisar M, et al. Structural
432 variants identified by Oxford Nanopore PromethION sequencing of the human genome. *Genome*
433 *Res* [Internet]. 2019; Available from: <http://dx.doi.org/10.1101/gr.244939.118>
- 434 26. Kim H-S, Jeon S, Kim C, Kim YK, Cho YS, Blazyte A, et al. Chromosome-scale assembly
435 comparison of the Korean Reference Genome KOREF from PromethION and PacBio with Hi-C
436 mapping information [Internet]. *bioRxiv*. 2019 [cited 2019 Jun 20]. p. 674804. Available from:
437 <https://www.biorxiv.org/content/10.1101/674804v1>
- 438 27. Tan MH, Austin CM, Hammer MP, Lee YP, Croft LJ, Gan HM. Finding Nemo: hybrid
439 assembly with Oxford Nanopore and Illumina reads greatly improves the clownfish (*Amphiprion*
440 *ocellaris*) genome assembly. *Gigascience* [Internet]. 2018;7:1–6. Available from:
441 <http://dx.doi.org/10.1093/gigascience/gix137>
- 442 28. Miller DE, Staber C, Zeitlinger J, Hawley RS. Highly Contiguous Genome Assemblies of 15
443 *Drosophila* Species Generated Using Nanopore Sequencing. *G3* [Internet]. 2018;8:3131–41.
444 Available from: <http://dx.doi.org/10.1534/g3.118.200160>
- 445 29. Cao MD, Nguyen SH, Ganesamoorthy D, Elliott AG, Cooper M, Coin LJM. Scaffolding and
446 Completing Genome Assemblies in Real-time with Nanopore Sequencing [Internet]. *bioRxiv*.
447 2016 [cited 2019 Jun 20]. p. 054783. Available from:
448 <https://www.biorxiv.org/content/10.1101/054783v1>
- 449 30. Giordano F, Aigrain L, Quail MA, Coupland P, Bonfield JK, Davies RM, et al. De novo
450 yeast genome assemblies from MinION, PacBio and MiSeq platforms. *Sci Rep* [Internet].
451 2017;7:3935. Available from: <http://dx.doi.org/10.1038/s41598-017-03996-z>
- 452 31. Austin CM, Tan MH, Harrisson KA, Lee YP, Croft LJ, Sunnucks P, et al. De novo genome
453 assembly and annotation of Australia’s largest freshwater fish, the Murray cod (*Maccullochella*
454 *peelii*), from Illumina and Nanopore sequencing read. *Gigascience* [Internet]. 2017;6:1–6.
455 Available from: <http://dx.doi.org/10.1093/gigascience/gix063>
- 456 32. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly from
457 long uncorrected reads. *Genome Res* [Internet]. 2017;27:737–46. Available from:
458 <http://dx.doi.org/10.1101/gr.214270.116>
- 459 33. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an
460 integrated tool for comprehensive microbial variant detection and genome assembly
461 improvement. *PLoS One* [Internet]. 2014;9:e112963. Available from:
462 <http://dx.doi.org/10.1371/journal.pone.0112963>
- 463 34. Whitehead A. The evolutionary radiation of diverse osmotolerant physiologies in killifish
464 (*Fundulus* sp.). *Evolution* [Internet]. 2010;64:2070–85. Available from:
465 <http://dx.doi.org/10.1111/j.1558-5646.2010.00957.x>

- 466 35. Griffith RW. Environment and Salinity Tolerance in the Genus *Fundulus*. *Copeia* [Internet].
467 [American Society of Ichthyologists and Herpetologists (ASIH), Allen Press]; 1974;1974:319–
468 31. Available from: <http://www.jstor.org/stable/1442526>
- 469 36. Durack PJ, Wijffels SE, Matear RJ. Ocean salinities reveal strong global water cycle
470 intensification during 1950 to 2000. *Science* [Internet]. 2012;336:455–8. Available from:
471 <http://dx.doi.org/10.1126/science.1212222>
- 472 37. Burnett KG, Bain LJ, Baldwin WS, Callard GV, Cohen S, Di Giulio RT, et al. *Fundulus* as
473 the premier teleost model in environmental biology: opportunities for new insights using
474 genomics. *Comp Biochem Physiol Part D Genomics Proteomics* [Internet]. 2007;2:257–86.
475 Available from: <http://dx.doi.org/10.1016/j.cbd.2007.09.001>
- 476 38. Reid NM, Proestou DA, Clark BW, Warren WC, Colbourne JK, Shaw JR, et al. The genomic
477 landscape of rapid repeated evolutionary adaptation to toxic pollution in wild fish. *Science*
478 [Internet]. 2016;354:1305–8. Available from: <http://dx.doi.org/10.1126/science.aah4993>
- 479 39. Oziolor EM, Reid NM, Yair S, Lee KM, Guberman VerPloeg S, Bruns PC, et al. Adaptive
480 introgression enables evolutionary rescue from extreme environmental pollution. *Science*
481 [Internet]. 2019;364:455–7. Available from: <http://dx.doi.org/10.1126/science.aav4155>
- 482 40. Reid NM, Jackson CE, Gilbert D, Minx P, Montague MJ, Hampton TH, et al. The landscape
483 of extreme genomic variation in the highly adaptable Atlantic killifish. *Genome Biol Evol*
484 [Internet]. 2017; Available from: <http://dx.doi.org/10.1093/gbe/evx023>
- 485 41. Ghedotti MJ, Davis MP. Phylogeny, Classification, and Evolution of Salinity Tolerance of
486 the North American Topminnows and Killifishes, Family Fundulidae (Teleostei:
487 Cyprinodontiformes). *Fieldiana Life Earth Sci* [Internet]. 2013;7:1–65. Available from:
488 <http://www.bioone.org/doi/abs/10.3158/2158-5520-12.7.1>
- 489 42. Quick J. Ultra-long read sequencing protocol for RAD004 [Internet]. [protocols.io](https://www.protocols.io); 2018
490 [cited 2019 Jun 20]. Available from: <https://www.protocols.io/view/ultra-long-read-sequencing-protocol-for-rad004-mrxc57n>
- 492 43. Wasko AP, Martins C, Oliveira C, Foresti F. Non-destructive genetic sampling in fish. An
493 improved method for DNA extraction from fish fins and scales. *Hereditas* [Internet].
494 2003;138:161–5. Available from: <http://dx.doi.org/10.1034/j.1601-5223.2003.01503.x>
- 495 44. De Coster W, D’Hert S, Schultz DT, Cruts M, Van Broeckhoven C. NanoPack: visualizing
496 and processing long-read sequencing data. *Bioinformatics* [Internet]. 2018;34:2666–9. Available
497 from: <http://dx.doi.org/10.1093/bioinformatics/bty149>
- 498 45. Macmanes MD. On the optimal trimming of high-throughput mRNA sequence data. *Front*
499 *Genet* [Internet]. 2014;5:13. Available from: <http://dx.doi.org/10.3389/fgene.2014.00013>
- 500 46. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
501 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*
502 [Internet]. 2015;31:3210–2. Available from: <http://dx.doi.org/10.1093/bioinformatics/btv351>

- 503 47. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2 [Internet]. bioRxiv. 2019
504 [cited 2019 Jun 20]. p. 530972. Available from:
505 <https://www.biorxiv.org/content/10.1101/530972v1>
- 506 48. Olsen R-A, Bunikis I, Tiukova I, Holmberg K, Lötstedt B, Pettersson OV, et al. De novo
507 assembly of *Dekkera bruxellensis*: a multi technology approach using short and long-read
508 sequencing and optical mapping. Gigascience [Internet]. 2015;4:56. Available from:
509 <http://dx.doi.org/10.1186/s13742-015-0094-1>
- 510 49. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-molecule
511 sequencing and chromatin conformation capture enable de novo reference assembly of the
512 domestic goat genome. Nat Genet [Internet]. 2017;49:643–50. Available from:
513 <http://dx.doi.org/10.1038/ng.3802>
- 514 50. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al. Chromosome-scale
515 assemblies of plant genomes using nanopore long reads and optical maps. Nat Plants [Internet].
516 2018;4:879–87. Available from: <http://dx.doi.org/10.1038/s41477-018-0289-4>
- 517 51. Low WY, Tearle R, Bickhart DM, Rosen BD, Kingan SB, Swale T, et al. Chromosome-level
518 assembly of the water buffalo genome surpasses human and goat genomes in sequence
519 contiguity. Nat Commun [Internet]. 2019;10:260. Available from:
520 <http://dx.doi.org/10.1038/s41467-018-08260-0>
- 521 52. de Lannoy C, de Ridder D, Risse J. The long reads ahead: *de novo* genome assembly using
522 the MinION. F1000Res [Internet]. 2017;6:1083. Available from:
523 <http://dx.doi.org/10.12688/f1000research.12012.2>
- 524 53. Cali DS, Kim JS, Ghose S, Alkan C, Mutlu O. Nanopore Sequencing Technology and Tools
525 for Genome Assembly: Computational Analysis of the Current State, Bottlenecks and Future
526 Directions [Internet]. arXiv [q-bio.GN]. 2017. Available from: <http://arxiv.org/abs/1711.08774>
- 527 54. Nicholls SM, Quick JC, Tang S, Loman NJ. Ultra-deep, long-read nanopore sequencing of
528 mock microbial community standards [Internet]. bioRxiv. 2018 [cited 2019 Jun 20]. p. 487033.
529 Available from: <https://www.biorxiv.org/content/10.1101/487033v2>
- 530 55. De Bustos A, Cuadrado A, Jouve N. Sequencing of long stretches of repetitive DNA. Sci Rep
531 [Internet]. 2016;6:36665. Available from: <http://dx.doi.org/10.1038/srep36665>
- 532 55. De Bustos A, Cuadrado A, Jouve N. Sequencing of long stretches of repetitive DNA. Sci Rep
533 [Internet]. 2016;6:36665. Available from: <http://dx.doi.org/10.1038/srep36665>