

1 Article

## 2 Constructing a Reference Genome in a Single Lab: 3 The Possibility to Use Oxford Nanopore Technology

4 Yun Gyeong Lee<sup>1</sup>, Sang Chul Choi<sup>1</sup>, Yuna Kang<sup>1</sup>, Kyeong Min Kim<sup>1</sup>, Chon-Sik Kang<sup>2,\*</sup> and  
5 Changsoo Kim<sup>1,\*</sup>

6 <sup>1</sup> Department of Crop Science, College of Agricultural and Life Sciences, Chungnam National University,  
7 Daejeon, 34134, Republic of Korea

8 <sup>2</sup> National Institute of Crop Science, Rural Development Administration (RDA), Wanju, 55365, Republic of  
9 Korea

10 \* Correspondence: cskang@korea.kr (C.-S.K.); Tel.: +82-63-238-5227(C.-S.K.); changsookim@cnu.ac.kr (C.K.);  
11 Tel.: +82-42-821-5729 (C.K.)

12 **Abstract:** The whole genome sequencing (WGS) has become a crucial tool to understand genome  
13 structure and genetic variation. The MinION sequencing of Oxford Nanopore Technologies (ONT)  
14 is an excellent approach for performing WGS and has advantages in comparison with other Next-  
15 Generation Sequencing (NGS): It is relatively inexpensive, portable, has simple library preparation,  
16 can be monitored in real-time, and has no theoretical limits on read length. *Sorghum bicolor* (L.)  
17 Moench is diploid ( $2n = 2x = 20$ ) with a genome size of about 730 Mb, and its genome sequence  
18 information is released in the Phytozome database. Therefore, sorghum can be used as a good  
19 reference. However, plant species have complex and large genomes compared to animals or  
20 microorganisms. As a result, complete genome sequencing is difficult for plant species. MinION  
21 sequencing that produces long-reads can be an excellent tool to overcome the weak assembly of  
22 short-reads generated from NGS by minimizing the generation of gaps or covering the repetitive  
23 sequence that appears on the plant genome. Here, we conducted the genome sequencing for *S.*  
24 *bicolor* cv. BTx623 using the MinION platform and obtained 895,678 reads and 17.9 gigabytes (Gb)  
25 (ca. 25X coverage of reference) from long-read sequence data. Through a *de novo* assembly using  
26 two different tools and mapped assembled contigs against the sorghum reference genome, a total  
27 of 6,124 contigs (covering 45.9%) were generated from Canu, and a total of 2,661 contigs (covering  
28 50%) were generated from Minimap and Miniasm with a Racon pipeline. Our results provide a  
29 pipeline of long-read sequencing analysis for plant species using the MinION platform and a clue  
30 to determine the total sequencing scale for optimal coverage based on various genome sizes.

31 **Keywords:** sorghum; Canu; Miniasm; MinION; long-read sequencing

32

### 33 1. Introduction

34 The whole genome sequencing (WGS) has become a crucial tool to understand genome structure  
35 and genetic variation. Next-generation sequencing (NGS) technology, which has been actively used  
36 over the past decade, has revolutionized the genomic research of plants as well as animals and  
37 microorganisms, consequently leading to a high-throughput WGS [1,2]. However, most of the  
38 existing NGS techniques typically generate short-reads (35-700 bp) and the assembled sequences  
39 from these short-reads have resulted in an occurrence of gaps. This is because short-reads are not  
40 able to span repetitive sequences longer than their length due to the limitations of assembly  
41 completeness, thereby causing an incomplete genome assembly [1,3].

42 Unlike NGS, the third-generation sequencing (TGS) technology enables the generation of long-  
43 reads as a single molecule by preserving the native DNA state as much as possible during library  
44 construction and performing sequence detection through electrical or optical signals [2-4]. The major

45 advantage of long-read sequencing is that it may be able to resolve gaps that occurred from short-  
46 read assemblies [5]. Although the TGS market is overwhelmingly controlled by Pacific Biosciences  
47 (PacBio), the MinION platform [4] of Oxford Nanopore Technologies (ONT) is relatively inexpensive  
48 in comparison with other NGS platforms and allows the production of long-reads by only using small  
49 portable devices. In addition, the simple preparation of a sequencing library does not require a  
50 specific large instrumentation and complicated library preparation [2,6]. The MinION platform [4] is  
51 able to monitor sequence information in real-time as well as directly detect nucleotide modifications.  
52 As a result, this platform may be desirable for a small-scale laboratory to run and manage it in-house.  
53 Moreover, there is no theoretical limit on the read length, so if a high-molecular weight (HMW)  
54 genomic DNA (gDNA) and sequencing library were properly prepared, they could obtain a read  
55 sequence that has several hundred kilo base pairs or more [4,7]. If the high-error rate can be overcome,  
56 Nanopore sequencing may be very useful for the *de novo* assembly or for studying structural or single-  
57 nucleotide variations [4].

58 *Sorghum bicolor* (L.) Moench is one of the most consumed crops in the world and represents the  
59 C4 model plant. The WGS for sorghum has been performed and publicly available [8]. Sorghum is  
60 diploid ( $2n = 2x = 20$ ) with a genome size of about 730 Mb and a repeat content of ~61% (from  
61 homozygous sorghum genotype BTx623) [8,9]. Currently, the genome sequence information can be  
62 found in the Phytozome database (<https://www.phytozome.net/>, [10]) and is being continuously  
63 updated. However, to date, even though 4,426 gaps were closed, and the overall contiguity increased  
64 by 5.8X in a recent update (*S. bicolor* v3.1.1 in Phytozome database), Sorghum still remains an  
65 incomplete genome sequence. It is difficult to complete genome sequencing for plant species, since  
66 plant species have a more complex genome structure and larger genome size than animal species  
67 [11]. Plants have evolved through expanding or altering genomes, for example, the whole genome  
68 duplication, as a way to adapt to the external environment due to sessility, which results in a lot of  
69 repeated sequences [11-13]. During the evolutionary process, factors, such as polyploidy, repetitive  
70 sequences, heterozygosity, and transposable elements, have contributed to the plant genome size and  
71 complexity [11,14].

72 Recently, Nanopore sequencing using the MinION platform has been applied in various fields  
73 for plant species, but remains somewhat limited. The detection of transposable elements associated  
74 structural variants (TEASVs) in *Arabidopsis* [15], the validation of assemblers for the *Arabidopsis*  
75 genome [16], the *de novo* assembly of the *Solanum pennellii* genome through hybrid sequencing [17],  
76 the identification of novel genes related to nucleotide-binding leucine-rich repeat (NLR)[18], the  
77 improvement of maize reference genomes [19], and the field-based analysis for identifying closely-  
78 related plants (*Arabidopsis* spp.) [20] are some examples of this application. Moreover, apart from  
79 recent improvements in the accuracy of the Nanopore sequencing, there is a trend in an improved  
80 accuracy of assembled sequences by bioinformatically compensating long-reads using short-reads,  
81 leading to the obtaining of a high-contiguity genome assembly [21,22].

82 In this study, we conducted the genome sequencing for *S. bicolor* cv. BTx623 using the MinION  
83 platform [4] and obtained 895,678 in the read number and 17.9 Gb (ca. 25X coverage of the entire  
84 sorghum genome) from the long-read sequence data. We performed *de novo* assembly using two  
85 different tools and mapped the assembled contigs against the sorghum reference genome to  
86 determine how much the MinION sequencing results cover the entire genome. As a result, from  
87 Canu[23], a total of 6,124 contigs (344,453,188 bp in length covering 45.9% of reference) were  
88 generated, and from Minimap and Miniasm [24] with five rounds of Racon [25] polishing pipeline,  
89 a total of 2,661 contigs (375,105,174 bp in length covering 50% of reference) were generated. Our results  
90 provide a pipeline of long-read sequencing analysis for plant species using the MinION platform [4]  
91 and a clue to determine the total sequencing scale for optimal coverage based on various genome  
92 sizes in order to obtain satisfactory results for the *de novo* assembly.

## 93 2. Results

### 94 2.1. MinION sequencing of sorghum accession BTx623 genome

95 We conducted the sequencing of sorghum HMW gDNA by using the MinION platform [4] to  
 96 assess the high quality *de novo* assembly for the sorghum genome (cv. BTx623). The summary  
 97 statistics for each run were calculated separately and combined into one table (Table 1). We  
 98 constructed three libraries: DNA fragmentation was performed (around 20 Kbp) in one of the three  
 99 libraries (2<sup>nd</sup> in Table 1), and the remaining libraries used an intact HMW gDNA. MinION sequencing  
 100 for each library was conducted using the standard script provided in the MinKNOW software. The  
 101 total yielded amount of sequencing data varied between 2.85 Gb, 11.71 Gb, and 3.34 Gb with different  
 102 initial HMW gDNAs for library preparation. The 2<sup>nd</sup> result generated the largest data size compared  
 103 to the other two results (1<sup>st</sup> and 3<sup>rd</sup>) since it used fragmented HMW gDNA. A total of 17.9 Gb of raw  
 104 reads were generated, representing 25X of the total sorghum genome (based on 730 Mb). Overall, the  
 105 longest read length was up to 110 Kbp, while the most abundant reads were in the range of 908 bp to  
 106 1,028 bp in length.

107 The raw sequences were aligned to the sorghum BTx623 reference genome using the BWA-mem  
 108 version 0.7.15 [26] with the default option. All of the raw reads were separately analyzed and  
 109 combined before downstream processing. The average depth was approximately 8.6X and the  
 110 mapping rate was 97% for the combined data (Table 2). The Q-score was around 10, indicating that a  
 111 read error rate should be around 10%. The coverage distribution was plotted using the Mosdepth  
 112 [27] output (Figure 1). The depth of coverage calculation results from both the SAMtools [28] and  
 113 Mosdepth [27] showed that only 7.0X to 8.56X of the sorghum genome were covered by using the  
 114 combined sequencing data (17.9 Gb in total) generated from the three libraries. However, the  
 115 mapping rates were more than 97%, indicating that the sequencing data generated from the MinION  
 116 platform could contain redundant coverage in the specific regions of the sorghum genome. This  
 117 redundant coverage was particularly concentrated in regions having highly repeated DNA contents.  
 118 Furthermore, in many cases, certain regions are difficult to sequence and/or map because of repetitive  
 119 DNA or sequences aligned to multiple places in the genome. We will need additional data to resolve  
 120 these problems.

121 **Table 1.** The statistics of the raw fastq file

Result	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>
Total generated file size (Gb)	2.83	11.71	3.34
Total number of .fastq files	35	170	37
Total read numbers	136,769	679,658	146,883
The shortest read length (bp)	167	74	38
The longest read length (bp)	190,250	110,486	217,000
The most abundant read length (bp) (no. of reads)	908 (61)	947 (111)	1,028 (69)
Q-score	11.2	10.7	10.9

122

123 **Table 2.** Results of average depth and mapping rate for raw reads against reference genome.

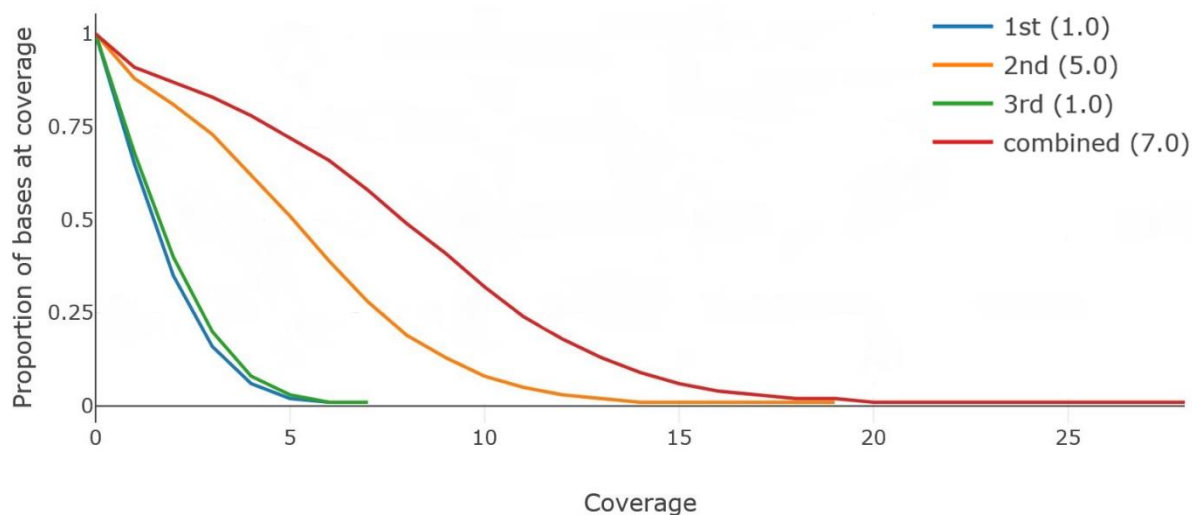
Result	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	Combined <sup>a</sup>
Average depth	2.01	5.64	2.10	8.56
Mapping rate (%)	97.93	96.87	97.14	97.08

124

<sup>a</sup>Combined all three results

125

126



127  
128 **Figure 1.** The coverage graph using Mosdepth. In this graph, the legend indicates the coverage graph  
129 for each result. The numbers in the parentheses indicate an average depth of coverage.  
130

### 131 2.2. Assembly results using Canu

132 The processed raw reads (from UniTigging/READs) were *de novo* assembled using Canu (version  
133 1.6) [23]. The correction step in the Canu assembly [23] improves the accuracy of each read base by  
134 building read and overlapped databases and choosing overlaps for correction. For the AT/GC rich  
135 eukaryotic genome, the `corMaxEvidenceErate = 0.15` parameter is suggested from the Canu  
136 documentation. Therefore, this parameter was added to run our data analysis and other options were  
137 used as the default.

138 A total of 9.4 Gb out of 17.9 Gb raw reads were loaded due to the specific feature that the low  
139 coverage data less than 10X in any region are eliminated by the Canu [23] program. After the  
140 correction step, only 8.0 Gb (11.56X) remained (Table 3). The correction phase improved the accuracy  
141 of bases in the reads, while the trimming phase cleaned the reads to the portion that appeared to be  
142 a high-quality sequence, as well as removed suspicious regions such as the remaining SMRTbell  
143 adapter. However, this was only applicable to the PacBio data. Therefore, the trimming phase may  
144 not drastically affect the entire read contents for MinION [29] trials. The final assembly phase ordered  
145 the reads into contigs and generated consensus sequences (unitigging consensus sequence). The final  
146 unitigging consensus sequence length was about 344 Mbp (344,366,012 bp) with N50, 98 Kbp (97,987  
147 bp).  
148

149 **Table 3.** Summary of read data for the results of Canu.

Result		1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	Combi ned
Total loaded reads	No. of reads	2	3	3	8
	Total length (bp)	1,495,9 87,647	6,216,3 12,936	1,767,1 14,081	9,479,4 14,664
	Coverage	2.04	8.51	2.42	12.63
	Expected corrected reads	No. of reads	2	1	6
Expected corrected reads	Total length (bp)	1,333,1 02,902	6,187,5 51,664	1,359,0 65,630	8,029,1 84,425
	Mean read length (bp)	11,304	7,900	10,800	8,986

	N50 length (bp)	49,358	23,337	53,805	72,703
After correction /Before trimming	No. of reads	0	5	0	4
	Total length (bp)	1,235,198,760	5,658,532,542	1,549,842,529	8,673,782,926
	Coverage	1.68	7.75	2.12	11.56
After trimming <sup>a</sup>	No. of reads	68,176	5	403,75	56,719
	Total bases (bp)	411,454,770	2,794,594,634	376,245,841	4,739,533,665
UniTigging/RE ADs	No. of reads	68,176	6	410,74	56,719
	Total length (bp)	424,463,809	2,844,670,276	381,679,172	4,833,385,452
	Coverage	0.58	3.89	0.52	6.44
UniTigging/con census	No. of sequences	159	5,740	127	6,124
	No. of repeats	28	692	26	712
	Length of repeats (bp)	573,105	10,509,344	472,695	14,815,759
	Total length (bp)	3,088,777	178,246,454	3,256,717	344,366,012
	Coverage	0.004	0.237	0.004	0.459
Unassembled	No. of sequences	38,897	8	168,88	32,340
	Total length (bp)	259,436,098	1,180,881,063	252,418,869	1,832,920,246

<sup>a</sup>Trimmed reads output

150  
151

### 152 2.3. Assembly results using Minimap, Miniasm, and Racon

153 The raw read overlapper, Minimap [24], was used to find overlaps, and Miniasm [24] was used  
154 to complete *de novo* assemblies using Minimap [24] results (Table 4). They directly produce  
155 unpolished and uncorrected contig sequences from the overlaps of raw reads. As a result, polishing  
156 steps should be indispensable to improve their credibility. The five rounds of Racon [30] were used  
157 to correct raw contigs to produce better quality sequences. The file size differences between the raw  
158 file (17.9 Gb) and Minimap (13.2 Gb) indicated that our combined data had about 4.7 Gb file size of  
159 duplicated overlaps. By using 13.2 Gb size of raw read overlaps, the unpolished and uncorrected  
160 contig sequences with a file size of 368 Mb and a contig length of 370 Mbp were generated. The final  
161 length of consensus sequences for the three combined data sets after five rounds of Racon [30]  
162 polishing steps was about 375 Mbp with a N50 value of 199 Kbp (Table 4). In this consensus sequence,  
163 the longest contig was 1 Mbp in length and the shortest contig was 779 bp. The final sequence length  
164 for the combined data from the Racon result (375,105,174 bp) (Table 4) was slightly longer than that  
165 of the Canu [23] result (344,366,012 bp) (Table 3).

166  
167

**Table 4.** Summary of Miniasm assemblies with Minimap and Racon.

Result	N o. of round	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	Combine d
--------	---------------------	-----------------	-----------------	-----------------	--------------

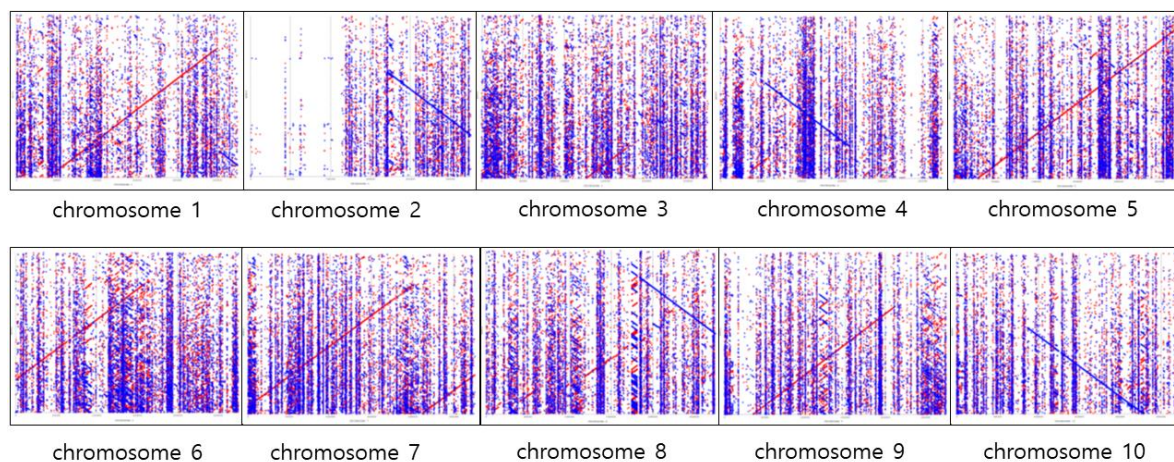


Ra w file	Total (Gb)	size	2.83	11.71	3.34	17.9
Mi nimap	File (byte)	size	607,227, 298	5,089,82 4,937	546,909, 744	13,226,11 0,131
Mi niasm	File (byte)	size	1,282,82 2	177,933, 354	2,126,52 5	368,271,9 34
	Total length (bp)		1,286,78 2	176,978, 175	2,139,68 2	370,303,4 49
		1	1,289,49	177,650,	2,145,74	373,675,1
		2		167	9	34
		2	1,278,46	177,931,	2,141,27	374,668,3
		7		139	7	65
Rac on	Total length (bp)	3	1,262,94	177,915,	2,127,08	374,934,5
		7		228	9	32
		4	1,247,13	177,805,	2,112,23	375,048,7
		8		838	9	32
		5	1,232,80	177,683,	2,097,34	375,105,1
		8		528	1	74

168

## 169 2.4. Confirmation of de novo assembly

170 To visualize the alignments between the assembled sequences from Miniasm [24] and each of  
 171 the sorghum reference chromosomes, the mummerplot option from the Mummer software version  
 172 3.0, [31] with default parameters (Figure 2) was used. A total of 2,661 contigs from 375,105,174 bp in  
 173 length were generated after doing polishing steps five times with Racon [30]. The  $x$ -axis represents  
 174 each chromosome of reference and the  $y$ -axis represents 2,261 contigs. A perfect alignment between  
 175 the contigs and each chromosome would completely fill the positive diagonal (slope == 1), while a  
 176 line of slope == -1 represents an inverted segment of conservation between the two sequences. In  
 177 chromosome 2, the contigs were not aligned either forward or reverse in some parts of the  
 178 chromosome. In other chromosomes, the contigs tended to partially align to chromosomes either  
 179 forward or reverse (Figure 2). It was almost impossible to confirm the alignment trend between the  
 180 Canu [23] consensus sequences and sorghum reference, since the Canu [23] generated relatively short  
 181 and almost three times more contigs (6,124 contigs) than Racon (2,661 contigs) [30] that used the  
 182 mummerplot (refer to Table 5 for comparison).  
 183



184

185

186

**Figure 2.** The five rounds polishing with Racon after Miniasm assembly versus each chromosome of the sorghum reference. The  $x$ -axis represents each chromosome of reference, and the  $y$ -axis represents

187 2,261 contigs. The forward matches are displayed in red, while the reverse matches are displayed in  
 188 blue.

189

190 **Table 5.** Comparison between Canu and Miniasm using final assembly results.

	Canu	Miniasm
Number of Conigs	6,124	2,661
Assembled read length (bp)	344,366,012	375,105,174
N50 (bp)	98,000,000	199,000,000

191

### 192 3. Discussion

#### 193 Optimization of genome assembly by using different assemblers

194 We performed MinION sequencing [29] for the sorghum and compared the final results obtained  
 195 from two different *de novo* assemblies against the reference genome in this study. With the ca. 12.63X  
 196 coverage raw reads, the *de novo* assemblies from Canu [23] and Miniasm [24] showed 0.459X and 0.5X  
 197 coverages for the sorghum genome, respectively. In other words, the completion of the *de novo*  
 198 assembly can be mathematically achieved by increasing the amount of raw data to ca. 25X for the  
 199 sorghum example. However, the quality of the *de novo* assembly can be affected by a plethora of  
 200 factors such as the contents of genes, GC ratio, and the length of repetitive sequences, genome size,  
 201 and ploidy numbers. We once tried to formulate the relationship between the minimum coverage  
 202 required for the *de novo* assembly and the amount of raw reads. However, it was not feasible due to  
 203 various factors. Nonetheless, our results showed that the *de novo* assembly for any species that does  
 204 not have reference sequences can be performed in the single laboratory with cost-effective ways due  
 205 to the newly developed ONT apparatus. In addition, a comparison between the two representative  
 206 long-reads assemblers, Canu [23] and Miniasm [24], indicates that Miniasm [24] with Racon [30]  
 207 correcting steps provides a better assembly in terms of the number of contigs and N50 values.

208 Depending on the type of assembler, different assembly results may be obtained for the  
 209 genome of the plant species [16,32]. As aforementioned, we showed the *de novo* assembly results using  
 210 Canu and Miniasm (with Minimap and Racon) for the sorghum genome. However, considering the  
 211 genome size, structural variation, and genome complexity of the genome in the genome assembly of  
 212 a particular plant species, it is imperative to determine which assembler is suitable for optimal results.  
 213 In addition, bioinformatic efforts should be used to ensure that the misassembled or ambiguous  
 214 sequences, such as repetitive regions of the genome, gaps, discrimination between paralogues and  
 215 alleles or between genes and pseudogenes [11], are properly assembled. RNA sequencing using the  
 216 MinION platform can be used to identify the isoforms of transcripts. Through this, various isoforms  
 217 have been identified without any assembly process, followed by non-redundant isoform clustering.  
 218 The resulting information can be used for genome annotation, and, consequently, can be integrated  
 219 to increase the contiguity and accuracy of the results from the existing genome assembly [5].

220

#### 221 Advantages of current combinational sequencing

222 Compared with studies on microorganisms or animal species, there is still not a lot of research  
 223 in plant species because of its genome complexity. Genome sequencing has continued to develop  
 224 through the classical Sanger method and NGS to TGS. Plant species have been actively studied in this  
 225 process, but already-sequenced plant species have genomes with low complexity and are of relatively  
 226 small size until the advent of NGS [11]. NGS technology makes it possible to perform sequencing  
 227 regardless of genome size, which is a substantial technical breakthrough to overcome these  
 228 limitations. However, sequencing for plant species with large and complex genomes was not resolved  
 229 in terms of contiguity and accuracy due to the limitations of the short-read assembly. In this respect,  
 230 TGS that produces long-reads can be an excellent tool. For instance, wheat (*Triticum aestivum*) has a  
 231 genome size of about 15Gb, an allohexaploid ( $2n = 6x = 42$ ), and a high repetitive character. The  
 232 International Wheat Genome Sequencing Consortium (IWGSC) has carried out wheat genome

233 assembly through a chromosome-based approach to conquer the genome nature of wheat, but it  
234 contained only 10.2 Gb of genomes with low contiguity [33,34]. Despite these efforts, the near-  
235 complete assembly of wheat has been achieved in a recent study [22].

236 This study demonstrates the possibility of assembling high complex genomes through a  
237 combination of sequencing Illumina short-reads and PacBio long-reads. The production of long-reads  
238 using TGS is able to overcome the weakness of assembling short-reads by minimizing the generation  
239 of gaps or covering the repetitive sequence that appears on the plant genome. In another aspect, when  
240 considering only the accuracy, short-reads can be used for error-correction by aligning them to long-  
241 reads, which enable the increased accuracy of the genome assembly [35]. Therefore, a hybrid  
242 assembly through combinational sequencing is a useful approach, at least until now, to overcome the  
243 limitations of the current two techniques. As a result, more accurate sequence data would be  
244 obtained. Even though the PacBio Single Molecule and Real-Time (SMRT) sequencing played a  
245 leading role, given the ease of performance and utilization of the MinION sequencing [29], the  
246 MinION sequencing is expected to replace the PacBio sequencing in laboratory-level sequencing. In  
247 the future, MinION sequencing [29] will play a significant role in noticeably improving the assembly  
248 of high complex genomes.

249

### 250 **Improvements in the accuracy of long-reads sequencing and assembly**

251 We did not perform the Illumina sequencing in this study since the sorghum genome sequence  
252 was already released. However, for a reference-free species, the *de novo* genome assembly is required.  
253 Therefore, the hybrid assembly will be sufficient to overcome the incompleteness caused by using a  
254 single platform. However, the hybrid assembly is more difficult than an assembly that uses a single  
255 platform. If the accuracy of the raw long-reads is high, or it can be increased by using the MinION  
256 platform [29] alone, the genome assembly of plant species with a highly complex genome is possible.  
257 However, the accuracy of Nanopore sequencing (85% accuracy for R9 version) is not high compared  
258 to that of the NGS generating short-reads [36]. Currently, even with the R9.5 flow cell using 1D<sup>2</sup>  
259 chemistry for the MinION [29], the model accuracy of sequences obtained using Nanopore  
260 sequencing is about 97% (<http://nanoporetech.com>) [29]. In contrast, the short-read from the Illumina  
261 platform has a maximum length of 150-300 bp, but most bases have more than 30 in quality score  
262 (99.9% accuracy) for single and paired-end reads (<https://www.illumina.com>). In addition, the  
263 improvement in sequencing accuracy can lead to the conclusion that the consensus accuracy will gain  
264 a high value from a small amount of raw read coverage [4]. For the *de novo* genome assembly, a raw  
265 read coverage of about 50-60X is needed to generate enough coverage of reads to cover repetitive  
266 regions in the genome assembly [37]. At this time, Nanopore sequencing for raw reads is not able to  
267 be more accurate than the accuracy of NGS, such as the Illumina sequencing, which produces highly  
268 accurate reads. Thus, an error-correction process is indispensable to increase the accuracy in  
269 Nanopore sequencing. Because of this, if only MinION is used as a single platform, the significance  
270 of correction tools for raw reads is greater. Nanocorrect (<https://github.com/jts/nanocorrect/>), PoreSeq  
271 [38], and nanoCORR (Goodwin et al., 2005) have been developed as a representative error-correction  
272 tool for the Nanopore sequence data. Recently, Canu [23], Falcon  
273 (<https://github.com/PacificBiosciences/FALCON/>), and Miniasm [24] assemblers are more commonly  
274 used for error correction as well as the assembly. However, we should be aware that when  
275 sequencing using MinION [29] as a single platform, it is advantageous to obtain long-reads using  
276 HMW gDNA, since adding reads to reduce the average length of reads is able to reduce the  
277 assembly's quality [39].

278

### 279 **DNA fragmentation effect on MinION sequencing**

280 The MinION flow cell (R9.4) that consists of 512 channels and four wells (four nanopores) are  
281 included in each channel. However, the read data are only generated from one of the four wells at a  
282 time [37]. From our results, the 2<sup>nd</sup> result that used fragmented gDNA produced more read data than  
283 the 1<sup>st</sup> and 3<sup>rd</sup> results that used intact HMW gDNA. However, the read length showed the opposite  
284 pattern. This may be due to the feature that the Nanopore sequencing could not be performed



285 simultaneously in four wells in each channel. As a result, the following possibilities can be  
286 considered. First, in the process of tethering DNA molecules onto a membrane near a pore protein,  
287 HMW DNA molecules may cause the spatial hindrance to deteriorate the accessibility of the other  
288 DNA molecules to the nanopore. Second, the time required for the HMW gDNA molecule to pass  
289 through the nanopore is too long to allow for sequencing in the other wells of the same channel. This  
290 may result in a decrease in sequencing efficiency. For example, when using R9 chemistry (about 250  
291 bp sequencing speed per second; <https://nanoporetech.com>), it takes about 1,000 seconds for 250 kb  
292 of the DNA molecule to sequence. In this case, assuming that it shows 100% efficiency, the number  
293 of reads obtained through MinION [29] sequencing for 48 h is only 172.8 reads in each channel. For  
294 now, depending on the experimental purpose, we need to choose whether to get a relatively large  
295 number of reads or to get reads that are as long as possible. We expect to meet both through future  
296 technical advances.

297

### 298 **Requirement of effective size selection for long-reads sequencing**

299 It is important to remove short-reads for high quality assembly. In this study, a large amount of  
300 short-reads (around 1 kb) was generated by MinION sequencing [29]. As aforementioned, if short-  
301 reads less than the average length are produced, the assembly quality can decrease. Thus, we should  
302 consider the possibility of generating a lot of short-reads, even though HMW gDNA is used as an  
303 initial material. In general, a certain level of DNA supercoiling is maintained *in vivo* [40]. However,  
304 during DNA extraction, the DNA may be damaged, and the DNA supercoil level may decrease. After  
305 DNA extraction, DNA repair and adapter ligation steps are performed during the DNA library  
306 preparation for MinION sequencing [29]. At this time, the efficiency of library production may vary  
307 depending on the structural complexity of the DNA. Highly ordered structures of genomic DNA may  
308 reduce the accessibility of enzymes involved in the DNA repair or adapter ligation, while short DNA  
309 fragments are expected to increase the efficiency of library production due to the relatively high  
310 accessibility of the enzymes. We also cannot rule out the possibility of DNA shearing by physical or  
311 chemical reactions during the DNA library preparation.

312 Another possibility is limiting the use of magnetic beads in the size selection and purification of  
313 the DNA library. Magnetic beads make it easy to remove small DNA that are less than 500 bp, but  
314 they are not effective in removing large size DNA. In addition, when a relatively small amount of  
315 beads is used to obtain large-sized DNA fragments (e.g., DNA fragments of 1-10 kb in size), the yield  
316 of the DNA itself is greatly reduced. It is possible that the limitations of the protocols used in this  
317 study may not have effectively removed small size DNA. This can be overcome by conventional size  
318 selection methods such as using gel electrophoresis and gel elution or automated DNA size selection  
319 (e.g. Pippin). However, until now, automated size selection is the most effective method, although it  
320 does not completely remove the short-reads. If a more convenient and efficient size selection method  
321 is developed, more accurate Nanopore sequencing and subsequent analysis will be possible.

## 322 **4. Materials and Methods**

### 323 **Plant material and genomic DNA extraction**

324 The sorghum reference accession BTx623 was obtained from the National Agrobiodiversity  
325 Center of the Rural Development Administration in Korea. Sorghum plants were grown on a  
326 Murashige and Skoog (MS) medium (Duchefa) in an artificial growth chamber (25°C, 14 h light/10 h  
327 dark) for 7-10 days. Shoot parts were only used for genomic DNA (gDNA) extraction, and the  
328 procedure of the gDNA extraction was performed following the method previously described ([41];  
329 [42]) with some modifications. Shoots of sorghum seedling were ground into a fine powder in liquid  
330 nitrogen by using a mortar and pestle. 100 mg of the sample powder was transferred into a 2 mL tube  
331 (eppendorf) containing 600  $\mu$ L of a modified Carlson buffer [100 mM Tris-HCl, pH 8.0, 2% CTAB, 1.4  
332 M NaCl, 1% PEG 8000, 20 mM EDTA, 2% PVP40, 0.1% ascorbic acid] pre-warmed to 60°C and 20  $\mu$ L  
333 of RNase A (20 mg/mL; invitrogen). The sample was immediately homogenized by inverting it gently  
334 20 times and incubating it in a water-bath at 60°C for 30 min with gentle inverting 20 times every 10  
335 min. After incubation, the sample was cooled-down to room temperature, and 600  $\mu$ L of chloroform

336 was added. The sample was inverted carefully 60 times. Afterwards, the sample was centrifuged at  
337 5,000 g for 10 min at 4°C, and 400  $\mu$ L of the supernatant was transferred to a new 2 mL tube. 400  $\mu$ L  
338 of the binding buffer (20% PEG 8000, 3 M NaCl) and 50  $\mu$ L of the AMPure XP beads solution were  
339 added to the sample and incubated with rotation (6 rpm) at room temperature for 10 min. The sample  
340 was briefly centrifuged and kept on a magnetic rack (Thermo Scientific) until the magnetic beads  
341 were completely separated. The supernatant was removed without disturbing the pellet, and then 1  
342 mL of 70% ethanol was added to the pellet. The pellet was incubated in ethanol for 1 min, and the  
343 supernatant was removed. The ethanol washing step was repeated 3 times. After the ethanol was  
344 removed, the sample was air-dried for 1 min. The pellet was eluted using a Buffer EB (Qiagen). At  
345 that moment, the amount of the Buffer EB was adjusted so that the eluate concentration was 80 ng/ $\mu$ L  
346 or more. As a result, HMW gDNA with a size longer than at least 50 kb was obtained.

347

### 348 **Preparation of sequencing library and MinION sequencing**

349 12  $\mu$ g of HMW gDNA was fragmented using a g-TUBE (Covaris) by centrifuging at 3,170 g for  
350 60 sec (Labogene 1730R; rotor GRF-M-m2.0-24). Of the three flow cells, the HMW gDNA of one flow  
351 cell was only fragmented, and the rest was used in its native state. The DNA library was prepared  
352 with the ONT Ligation Sequencing Kit 1D (SQK-LSK108), and the DNA preparation method was  
353 based on "1D gDNA long reads without BluePippin protocol" provided by the Nanoporetech  
354 community ([https://community.nanoporetech.com/protocols/1d-gdna-without-](https://community.nanoporetech.com/protocols/1d-gdna-without-bluepippin/v/1/all_steps)  
355 [bluepippin/v/1/all\\_steps](https://community.nanoporetech.com/protocols/1d-gdna-without-bluepippin/v/1/all_steps)). A total of 2  $\mu$ g of gDNA (80 ng/ $\mu$ L) was used to construct the DNA library  
356 in each flow cell. MinION sequencing was performed using a R9.4 SpotON flow cell (FLO-MIN106),  
357 and the default script "NC\_48Hr\_Sequencing\_Run\_FLO-MIN106\_SQK-LSK108" from the  
358 MinKNOW program was used to run the sequencing. Finally, the read sequence files (fastq format)  
359 were obtained from the MinKNOW workflow.

360 Nanopore sequencing data have been deposited in NCBI (  
361 <https://www.ncbi.nlm.nih.gov/sra/PRJNA544582>) with the accession number of PRJNA544582.

362

### 363 **MinION raw sequences mapped against the reference genome**

364 The generated raw fastq files from the MinKNOW workflow were mapped to the sorghum  
365 BTx623 reference genome (v.3.1.1) and downloaded from the plant genomics resource  
366 (<https://phytozome.jgi.doe.gov/pz/portal.html>) by using the BWA mem (version 0.7.15) [26] with  
367 default parameters.

368 The average depth was evaluated with a depth option, and the mapping rate was conducted  
369 with the flagstat option in SAMtools version 1.3.1 [28]. The Mosdepth program version 0.2.3 [27] was  
370 used to calculate the depth from the BAM file at each nucleotide position in a genome and produce  
371 coverage graph (Figure 1).

372

### 373 **De novo whole genome assembly**

374 In order to handle these noisy long MinION reads efficiently, some specially designed tools were  
375 adopted. Two *de novo* assemblers were selected to compare their performances: Canu (version 1.6) [23]  
376 and Minimap and Miniasm (version 0.2-r168-dirty) [24]. Canu [23] is a new single-molecule sequence  
377 assembler that improves the Celera Assembler. Canu operates in three phases: Correction, trimming,  
378 and assembly. The correction step improves the accuracy of each read base. For the AT/GC rich  
379 eukaryotic genomes, the corMaxEvidenceErate = 0.15 parameter is suggested by the developer's  
380 instructions. Therefore, this parameter was incorporated to run our sorghum data, and other options  
381 were set as a default.

382 Minimap [24] with -Sw5 -L100 -m0 -t8 options and a *de novo* assembler, Miniasm [24] with default  
383 parameters were used to assemble MinION sequencing reads without an error correction stage.  
384 Minimap is an all-against-all read self-mapping tool, and Miniasm is composed of simple  
385 concatenated pieces of the read sequence to generate the final unitig sequences. This pipeline allows  
386 sequencing data to be assembled into a single contig in a relatively short time. However, the  
387 consensus sequence error rate is as high as the raw reads. Therefore, Racon

388 (<https://github.com/isovic/racon>)[30] coupled with the Miniasm pipeline could be used to generate  
389 similar or better quality final unitig sequences. Multiple rounds of Racon polishing have given a good  
390 final sequence accuracy and produced the best possible consensus sequence. To improve the  
391 sequence's quality, we conducted five rounds of Racon [30] using Minimap and Miniasm [24] results.

### 392 **Confirmation of *de novo* assembly against the sorghum reference genome**

393 To test the structural correctness of the unitig genome, we aligned the assembly results from the  
394 Miniasm against the sorghum BTx623 reference genome. The nucmer option from the MUMmer  
395 software version 3.0 [31] was used to get an overview of the global alignment between the contigs  
396 and reference genome. In addition, the delta-filtering option with the  $-r$  and  $-q$  parameters were used  
397 to filter the alignment results. The mummerplot option from the MUMmer [31] was used to draw the  
398 dotplot.

399 The MUMmer sequence alignment package [31] was designed to detect the homology regions  
400 in genome sequences. For a dotplot, the reference sequence is laid across the  $x$ -axis, while the query  
401 sequence is on the  $y$ -axis. Wherever the two sequences agree, a colored line or dot is plotted. The  
402 forward matches are displayed in red, while the reverse matches are displayed in blue.

403

## 404 **5. Conclusions**

405 Since the advent of TGS, Minion sequencing has developed rapidly in less than five years.  
406 Advances in its chemistry have elevated the speed and accuracy of sequencing, and the contiguity of  
407 genome assembly was improved by enabling long-reads. These developments have enhanced the  
408 high utilization and value of genome assemblies for plant species with highly complex genomes. We  
409 showed the results of MinION sequencing [29] for the *S. bicolor* cv. BTx623, in which the accuracy and  
410 coverage of raw data against the reference genome changed during the process of error-correction,  
411 *de novo* assembly, and polishing. Our results not only illustrate the use of appropriate tools for  
412 genome assembly through MinION sequencing [29] in plant species, but also provide information on  
413 the amount of raw data required for a more accurate genome assembly. This is expected to contribute  
414 to complete genome sequencing in a variety of plant species, including reference-free species.

415

## 416 **5. Patents**

417 **Acknowledgements:** This work is supported by the "Cooperative Research Program for Agriculture  
418 Science and Technology Development (Project No. PJ01252701)", Rural Development  
419 Administration, Republic of Korea

420

## 421 **References**

- 422 1. Goodwin, S.; McPherson, J.D.; McCombie, W.R. Coming of age: ten years of next-generation sequencing  
423 technologies. *Nature Reviews Genetics* **2016**, *17*, 333.
- 424 2. Levy, S.E.; Myers, R.M. Advancements in next-generation sequencing. *Annual review of genomics and*  
425 *human genetics* **2016**, *17*, 95-115.
- 426 3. Rhoads, A.; Au, K.F. PacBio sequencing and its applications. *Genomics, proteomics & bioinformatics* **2015**,  
427 *13*, 278-289.
- 428 4. Jain, M.; Olsen, H.E.; Paten, B.; Akeson, M. The Oxford Nanopore MinION: delivery of nanopore  
429 sequencing to the genomics community. *Genome biology* **2016**, *17*, 239.
- 430 5. Li, C.; Lin, F.; An, D.; Wang, W.; Huang, R. Genome sequencing and assembly by long reads in plants.  
431 *Genes* **2017**, *9*, 6.
- 432 6. Shendure, J.; Balasubramanian, S.; Church, G.M.; Gilbert, W.; Rogers, J.; Schloss, J.A.; Waterston, R.H.  
433 DNA sequencing at 40: past, present and future. *Nature* **2017**, *550*, 345.

- 434 7. Jain, M.; Koren, S.; Miga, K.H.; Quick, J.; Rand, A.C.; Sasani, T.A.; Tyson, J.R.; Beggs, A.D.; Dilthey, A.T.;  
435 Fiddes, I.T. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature*  
436 *biotechnology* **2018**, *36*, 338.
- 437 8. Paterson, A.H.; Bowers, J.E.; Bruggmann, R.; Dubchak, I.; Grimwood, J.; Gundlach, H.; Haberler, G.;  
438 Hellsten, U.; Mitros, T.; Poliakov, A. The Sorghum bicolor genome and the diversification of grasses.  
439 *Nature* **2009**, *457*, 551.
- 440 9. McCormick, R.F.; Truong, S.K.; Sreedasyam, A.; Jenkins, J.; Shu, S.; Sims, D.; Kennedy, M.;  
441 Amirebrahimi, M.; Weers, B.D.; McKinley, B. The Sorghum bicolor reference genome: improved  
442 assembly, gene annotations, a transcriptome atlas, and signatures of genome organization. *The Plant*  
443 *Journal* **2018**, *93*, 338-354.
- 444 10. Goodstein, D.M.; Shu, S.; Howson, R.; Neupane, R.; Hayes, R.D.; Fazo, J.; Mitros, T.; Dirks, W.; Hellsten,  
445 U.; Putnam, N. Phytozome: a comparative platform for green plant genomics. *Nucleic acids research* **2011**,  
446 *40*, D1178-D1186.
- 447 11. Claros, M.G.; Bautista, R.; Guerrero-Fernández, D.; Benzerki, H.; Seoane, P.; Fernández-Pozo, N. Why  
448 assembling plant genome sequences is so challenging. *Biology* **2012**, *1*, 439-459.
- 449 12. Crow, K.D.; Wagner, G.P. What is the role of genome duplication in the evolution of complexity and  
450 diversity? *Molecular biology and evolution* **2005**, *23*, 887-892.
- 451 13. Wendel, J.F.; Jackson, S.A.; Meyers, B.C.; Wing, R.A. Evolution of plant genome architecture. *Genome*  
452 *biology* **2016**, *17*, 37.
- 453 14. Jackson, S.A.; Iwata, A.; Lee, S.H.; Schmutz, J.; Shoemaker, R. Sequencing crop genomes: approaches  
454 and applications. *New Phytologist* **2011**, *191*, 915-925.
- 455 15. Debladis, E.; Llauro, C.; Carpentier, M.-C.; Mirouze, M.; Panaud, O. Detection of active transposable  
456 elements in Arabidopsis thaliana using Oxford Nanopore Sequencing technology. *BMC genomics* **2017**,  
457 *18*, 537.
- 458 16. Michael, T.P.; Jupe, F.; Bemm, F.; Motley, S.T.; Sandoval, J.P.; Lanz, C.; Loudet, O.; Weigel, D.; Ecker,  
459 J.R. High contiguity Arabidopsis thaliana genome assembly with a single nanopore flow cell. *Nature*  
460 *communications* **2018**, *9*, 541.
- 461 17. Schmidt, M.H.-W.; Vogel, A.; Denton, A.K.; Istace, B.; Wormit, A.; van de Geest, H.; Bolger, M.E.;  
462 Alseekh, S.; Maß, J.; Pfaff, C. De novo assembly of a new Solanum pennellii accession using nanopore  
463 sequencing. *The Plant Cell* **2017**, *29*, 2336-2348.
- 464 18. Giolai, M.; Paajanen, P.; Verweij, W.; Witek, K.; Jones, J.D.; Clark, M.D. Comparative analysis of  
465 targeted long read sequencing approaches for characterization of a plant's immune receptor repertoire.  
466 *BMC genomics* **2017**, *18*, 564.
- 467 19. Jiao, Y.; Peluso, P.; Shi, J.; Liang, T.; Stitzer, M.C.; Wang, B.; Campbell, M.S.; Stein, J.C.; Wei, X.; Chin,  
468 C.-S. Improved maize reference genome with single-molecule technologies. *Nature* **2017**, *546*, 524.
- 469 20. Parker, J.; Helmstetter, A.J.; Devey, D.; Wilkinson, T.; Papadopoulos, A.S. Field-based species  
470 identification of closely-related plants using real-time nanopore sequencing. *Scientific reports* **2017**, *7*,  
471 8345.
- 472 21. Yang, J.; Liu, D.; Wang, X.; Ji, C.; Cheng, F.; Liu, B.; Hu, Z.; Chen, S.; Pental, D.; Ju, Y. The genome  
473 sequence of allopolyploid Brassica juncea and analysis of differential homoeolog gene expression  
474 influencing selection. *Nature genetics* **2016**, *48*, 1225.
- 475 22. Zimin, A.V.; Puiu, D.; Hall, R.; Kingan, S.; Clavijo, B.J.; Salzberg, S.L. The first near-complete assembly  
476 of the hexaploid bread wheat genome, Triticum aestivum. *Gigascience* **2017**, *6*, gix097.



- 477 23. Koren, S.; Walenz, B.P.; Berlin, K.; Miller, J.R.; Bergman, N.H.; Phillippy, A.M. Canu: scalable and  
478 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* **2017**, *27*,  
479 722-736, doi:10.1101/gr.215087.116.
- 480 24. Li, H. Minimap and miniiasm: fast mapping and de novo assembly for noisy long sequences.  
481 *Bioinformatics* **2016**, *32*, 2103-2110, doi:10.1093/bioinformatics/btw152.
- 482 25. Vaser, R.; Sovic, I.; Nagarajan, N.; Sikic, M. Fast and accurate de novo genome assembly from long  
483 uncorrected reads. *Genome Res* **2017**, *27*, 737-746, doi:10.1101/gr.214270.116.
- 484 26. Li, H.; Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler transform.  
485 *Bioinformatics* **2010**, *26*, 589-595.
- 486 27. Pedersen, B.S.; Quinlan, A.R. Mosdepth: quick coverage calculation for genomes and exomes.  
487 *Bioinformatics* **2017**, *34*, 867-868.
- 488 28. Li, H.; Handsaker, B.; Wysoker, A.; Fennell, T.; Ruan, J.; Homer, N.; Marth, G.; Abecasis, G.; Durbin, R.  
489 The sequence alignment/map format and SAMtools. *Bioinformatics* **2009**, *25*, 2078-2079.
- 490 29. de Lannoy, C.; de Ridder, D.; Risse, J. The long reads ahead: de novo genome assembly using the  
491 MinION. *F1000Research* **2017**, *6*.
- 492 30. Vaser, R.; Sović, I.; Nagarajan, N.; Šikić, M. Fast and accurate de novo genome assembly from long  
493 uncorrected reads. *Genome research* **2017**, *27*, 737-746.
- 494 31. Kurtz, S.; Phillippy, A.; Delcher, A.L.; Smoot, M.; Shumway, M.; Antonescu, C.; Salzberg, S.L. Versatile  
495 and open software for comparing large genomes. *Genome biology* **2004**, *5*, R12.
- 496 32. Bouri, L.; Lavenier, D.; Gibrat, J.-F.; del Angel, V.F.D. Evaluation of genome assembly software based  
497 on long reads. France Genomique, 2017.
- 498 33. Gill, B.S.; Appels, R.; Botha-Oberholster, A.-M.; Buell, C.R.; Bennetzen, J.L.; Chalhoub, B.; Chumley, F.;  
499 Dvořák, J.; Iwanaga, M.; Keller, B. A workshop report on wheat genome sequencing: International  
500 Genome Research on Wheat Consortium. *Genetics* **2004**, *168*, 1087-1096.
- 501 34. International Wheat Genome Sequencing, C. A chromosome-based draft sequence of the hexaploid  
502 bread wheat (*Triticum aestivum*) genome. *Science* **2014**, *345*, 1251788, doi:10.1126/science.1251788.
- 503 35. Jayakumar, V.; Sakakibara, Y. Comprehensive evaluation of non-hybrid genome assembly tools for  
504 third-generation PacBio long-read sequence data. *Briefings in bioinformatics* **2017**.
- 505 36. Mahmoud, M.; Zywicki, M.; Twardowski, T.; Karlowski, W.M. Efficiency of PacBio long read correction  
506 by 2nd generation Illumina sequencing. *Genomics* **2017**.
- 507 37. Lu, H.; Giordano, F.; Ning, Z. Oxford Nanopore MinION Sequencing and Genome Assembly. *Genomics*  
508 *Proteomics Bioinformatics* **2016**, *14*, 265-279, doi:10.1016/j.gpb.2016.05.004.
- 509 38. Szalay, T.; Golovchenko, J.A. De novo sequencing and variant calling with nanopores using PoreSeq.  
510 *Nature biotechnology* **2015**, *33*, 1087.
- 511 39. Tyson, J.R.; O'Neil, N.J.; Jain, M.; Olsen, H.E.; Hieter, P.; Snutch, T.P. MinION-based long-read  
512 sequencing and assembly extends the *Caenorhabditis elegans* reference genome. *Genome research* **2018**,  
513 *28*, 266-274.
- 514 40. Corless, S.; Gilbert, N. Investigating DNA supercoiling in eukaryotic genomes. *Briefings in functional*  
515 *genomics* **2017**, *16*, 379-389.
- 516 41. Carlson, J.; Tulsieram, L.; Glaubitz, J.; Luk, V.; Kauffeldt, C.; Rutledge, R. Segregation of random  
517 amplified DNA markers in F1 progeny of conifers. *Theoretical and Applied Genetics* **1991**, *83*, 194-200.

- 518 42. Mayjonade, B.; Gouzy, J.; Donnadiou, C.; Pouilly, N.; Marande, W.; Callot, C.; Langlade, N.; Muños, S.  
519 Extraction of high-molecular-weight genomic DNA for long-read sequencing of single molecules.  
520 *BioTechniques* **2016**, *61*, 203-205.  
521