

de novo assembly and population genomic survey of natural yeast isolates with the Oxford Nanopore MinION sequencer

Benjamin Istace¹, Anne Friedrich², Léo d'Agata¹, Sébastien Faye¹, Emilie Payen¹, Odette Beluche¹, Claudia Caradec², Sabrina Davidas¹, Corinne Cruaud¹, Gianni Liti³, Arnaud Lemainque¹, Stefan Engelen¹, Patrick Wincker^{1,4,5}, Joseph Schacherer^{2,§}, Jean-Marc Aury^{1,§}

¹ Commissariat à l'Energie Atomique et aux Energies Alternatives (CEA), Institut de Génomique (IG), Genoscope, BP5706, 91057 Evry, France

² Université de Strasbourg, CNRS, GMGM UMR 7156, F-67000 Strasbourg, France

³ Institute of Research on Cancer and Ageing of Nice (IRCAN), CNRS UMR 7284-INSERM U1081, Faculté de Médecine, Université de Nice Sophia Antipolis, Nice, France.

⁴ Université d'Evry Val d'Essonne, UMR 8030, CP5706, 91057 Evry, France

⁵ Centre National de Recherche Scientifique (CNRS), UMR 8030, CP5706, 91057 Evry, France

§Corresponding authors

Benjamin Istace: bistace@genoscope.cns.fr

Anne Friedrich: anne.friedrich@unistra.fr

Léo d'Agata: ldagata@genoscope.cns.fr

Sébastien Faye: sfaye@genoscope.cns.fr

Emilie Payen: emathieu@genoscope.cns.fr

Odette Beluche: obeluche@genoscope.cns.fr

Claudia Caradec: claudia.caradec@unistra.fr

Sabrina Davidas: sdavidas@genoscope.cns.fr

Corinne Cruaud: ccruaud@genoscope.cns.fr

Gianni Liti: gianni.liti@unice.fr

Arnaud Lemainque: alemainque@genoscope.cns.fr

Stefan Engelen: sengelen@genoscope.cns.fr

Patrick Wincker: pwincker@genoscope.cns.fr

Joseph Schacherer: schacherer@unistra.fr

Jean-Marc Aury: jmaury@genoscope.cns.fr

36 **Abstract**

37 **Background:** Oxford Nanopore Technologies Ltd (Oxford, UK) have recently
38 commercialized MinION, a small single-molecule nanopore sequencer, that offers the
39 possibility of sequencing long DNA fragments from small genomes in a matter of seconds.
40 The Oxford Nanopore technology is truly disruptive, it has the potential to revolutionize
41 genomic applications due to its portability, low-cost, and ease of use compared with existing
42 long reads sequencing technologies. The MinION sequencer enables the rapid sequencing of
43 small eukaryotic genomes, such as the yeast genome. Combined with existing assembler
44 algorithms, near complete genome assemblies can be generated and comprehensive
45 population genomic analyses can be performed.

46 **Results:** Here, we resequenced the genome of the *Saccharomyces cerevisiae* S288C strain to
47 evaluate the performance of nanopore-only assemblers. Then we *de novo* sequenced and
48 assembled the genomes of 21 isolates representative of the *S. cerevisiae* genetic diversity
49 using the MinION platform. The contiguity of our assemblies was 14 times higher than the
50 Illumina-only assemblies and we obtained one or two long contigs for 65% of the
51 chromosomes. This high contiguity allowed us to accurately detect large structural variations
52 across the 21 studied genomes.

53 **Conclusion:** Because of the high completeness of the nanopore assemblies, we were able to
54 produce a complete cartography of transposable elements insertions and inspect structural
55 variants that are generally missed using a short-read sequencing strategy. Our analyses show
56 that the Oxford Nanopore technology is already usable for *de novo* sequencing and assembly;
57 however non-random errors in homopolymers require polishing the consensus using an
58 alternate sequencing technology.

59 **Keywords:** *de novo* assembly; Nanopore sequencing; Oxford Nanopore; MinION device;
60 genome finishing; structural variations; transposable elements

61 **Background**

1
2 62 Today, long-read sequencing technology offers interesting alternatives to solve genome
3
4 63 assembly difficulties and improve the completeness of genome assemblies, mostly in
5
6 64 repetitive regions [1] where short-read sequencing has failed. Microbial or small eukaryotic
7
8
9 65 genomes could now be assembled using Oxford Nanopore [2] or Pacific Biosciences reads
10
11 66 alone [3, 4] or in combination with short but high quality reads [5-7]. Application of the
12
13 67 single-molecule real-time (SMRT) sequencing platform to large complex eukaryotic genomes
14
15 68 demonstrated the possibility of considerably improving genome assembly quality [8, 9].
16
17
18 69 Similar improvements were also accomplished using the 10x Genomics platform, and its
19
20 70 application to the human genome produced encouraging results [10-12] and showed the
21
22 71 importance of obtaining long and high-quality reads.
23
24
25

26
27 72 The most used sequencing technologies are based on the synthesis of new DNA strands,
28
29 73 including the Illumina and Pacific Biosciences technologies [13]. These sequencing
30
31 74 technologies based on optical detection of nucleotide incorporations are often commercialized
32
33 75 through large-sized and expensive instruments. For example, the cost of the commercially
34
35 76 available Pacific Biosystems RS II instrument is high and the infrastructure and
36
37 77 implementation needs make it inaccessible to large sections of the research community. This
38
39 78 year Oxford Nanopore Technologies Ltd (ONT, Oxford, UK) commercialized MinION, a
40
41 79 single-molecule nanopore sequencer that can be connected to a laptop through a USB
42
43 80 interface [14, 15]. This system is portable (close to the size of a harmonica) and low-cost
44
45 81 (currently USD 1,000 for the instrument). The MinION technology is based on an array of
46
47 82 nanopores embedded on a chip that detects consecutive 6-mers of a single-strand DNA
48
49 83 molecule by electrical sensing [16-19]. In addition to its small size and low price, this new
50
51 84 technology has several advantages over the older technologies. Library construction involves
52
53 85 a simplified method, no amplification step is needed, and data acquisition and analyses occur
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

86 in real time [20]. Library preparation can be performed in two ways: (i) a 10-minute library
87 preparation based on an enzymatic method for ‘1D’ sequencing (sequencing one strand of the
88 DNA) or (ii) a library preparation based on ligation for ‘2D’ sequencing (sequencing both the
89 template and complement strands of the DNA). In the 2D sequencing mode, the two strands
90 of a DNA molecule are linked by a hairpin and sequenced consecutively. When the two
91 strands of the molecule are read successfully, a consensus sequence is built to obtain a more
92 accurate read (called 2D read). Otherwise only the template or complement strand sequence is
93 provided (called 1D read).

94 Here, we sequenced the genomes of 22 *Saccharomyces cerevisiae* isolates to determine if the
95 MinION system could be used in population genomic projects that require a deeper view of
96 the genetic variation landscape. Even if the throughput of MinION was still heterogeneous,
97 we were able to perform the sequencing in a reasonable time using six MinION devices (less
98 than 2 days per strain). First, we resequenced the *Saccharomyces cerevisiae* S288C reference
99 genome using a nanopore long-read sequencing strategy to evaluate recent assembly methods.
100 We generated a complete benchmark of the assembly structures, as well as the completeness
101 of complex regions. Next, we selected 21 strains of *S. cerevisiae* that were genetically diverse,
102 based on preliminary results of the 1002 Yeast Genomes Project a large-scale short-read
103 resequencing project (<http://1002genomes.u-strasbg.fr/>). The genomes of these 21 strains
104 were *de novo* sequenced and assembled with Nanopore long-reads to have a better insight into
105 the variation of their genomic architecture. We obtained near complete assembly, in terms of
106 genes, as well as transposable elements and telomeric regions. The most contiguous assembly
107 produced a single contig per chromosome, except for chromosomes 3 and 12, the latter
108 containing the large repeated rDNA cluster.

110 Results

111 MinION data evaluation

112 We first sequenced the S288C genome by doing 11 MinION Mk1 runs with the R7.3
113 chemistry. On average, a 48-hours run produced more than 200 Mb of sequence, and the best
114 run throughput was 400 Mb. Two 2D library types with 8 kb and 20 kb mean fragmentation
115 sizes were used. They led to nearly 360,000 reads with a cumulative length of approximately
116 2.3 Gb and 63% of the nucleotides were in 2D reads, which represented a 187x and 118x
117 genome coverage for 1D and 2D reads, respectively. Template reads had a median length of
118 8.9 kb while 2D reads had a median length of 7.7 kb. All sequencing reads were aligned to the
119 S288C reference genome using BWA [21] to assess their quality. We successfully aligned
120 95.6% of the 2D reads with an average error rate of 17.2% (**Figure 1a**). ONT tagged high-
121 quality 2D reads as “2D pass” reads (reads with an average per-base quality higher than 9),
122 and 99.7% of the 2D pass reads were aligned to the reference genome with an average error
123 rate of 12.2%. We then parsed the alignment files to search for errors in stretches of the same
124 nucleotide (homopolymers). About 85% of A, T, C, and G homopolymers of size 2 were
125 present correctly in the reads. This percentage decreased rapidly to 65% for homopolymers of
126 size 4 for A and T homopolymers and to 70% for C and G homopolymers. For size 7
127 homopolymers, it was 30% for A and T homopolymers and 35% for C and G homopolymers
128 (**Figure S1a**).

129 We also sequenced the S288C genome using the R9 chemistry, the recently released version
130 of the pore. We obtained approximately 1 Gb of reads; 568 Mb were 2D reads, which
131 represents a 85x coverage with 1D reads and a 47x coverage with 2D reads. The mean 2D
132 length was 6.1 kb. We aligned 82.1% of the 1D reads with a mean identity percentage of
133 82.8% and 94.3% of the 2D reads with a mean identity percentage of 85.2% (**Figure 1b**). As
134 we did with the R7.3 reads, we also searched for errors in homopolymers (**Figure S1b**). The

135 numbers of correct A, T, C, and G homopolymers started at about 90% for size equal to 2,
136 then decreased to 75% for A and T homopolymers of size 4 and to 60% for the C and G
137 homopolymers. For size 7 homopolymers, it was 32% for A, T, and C homopolymers and
138 35% for G homopolymers.

139 **Comparison of Nanopore-only assemblers**

140 We tested Canu [22], Miniasm [23], SMARTdenovo [24] and ABruijn [25] with different
141 subset of 1D, 2D, and 2D pass reads (**Supplementary File 2 and Table S1**) and kept the
142 most contiguous assembly for each software.

143 With Canu, the assembly with the higher N50 was obtained with the whole set of 2D pass
144 reads (67x coverage). The assembly was composed of 37 contigs with a cumulative length of
145 12 Mb and seven chromosomes were assembled in one or two contigs. After aligning the
146 contigs to the S288C reference genome using Quast [26], we detected a high number of
147 deletions (120,365), which were often localized in homopolymers (58%). As a consequence,
148 only 454 of the 6,243 genes found in the assembly were insertion/deletion (indel)-free (**Table**
149 **S2**). With Miniasm, the most contiguous assembly was obtained using the 2D reads corrected
150 by Canu, which represented coverage of approximately 108x. The Miniasm assembly was
151 composed of 28 contigs with a cumulative length of 11.8 Mb, and 13 chromosomes were
152 assembled in one or two contigs. The Miniasm consensus sequence contained the higher
153 number of mismatches and indels (**Table S2**). With SMARTdenovo, 30x of the longest 2D
154 reads produced the assembly with the highest contiguity. It was composed of 26 contigs, with
155 a total length of 12 Mb, and 14 chromosomes were assembled in one or two contigs. The
156 SMARTdenovo assembly better covered the reference genome (>99%) and contained the
157 highest number of genes (98.8% of the 6,350 S288C genes), but the Quast output again
158 revealed a high number of deletions (128,050). With ABruijn, we obtained the assembly with
159 the highest N50 when using all the 2D reads as input, which represented coverage of

160 approximately 120x. The assembly contained 23 contigs with a cumulative length of 11.9 Mb,
161 and 14 chromosomes were assembled in one or two contigs (**Table S2**).

162 Next, we aligned the assemblies (Canu, Miniasm, SMARTdenovo, and ABruijn) to the S288C
163 reference genome using NUCmer [27], and visualized the alignments with mummerplot
164 (**Figures S2, S3, S4 and S5**). We also examined the coordinates of the alignments to search
165 for chimera. We did not detect any chimeric contigs in the Canu, Miniasm, or SMARTdenovo
166 assemblies; however, we did find some in the ABruijn assembly. Three chimeric contigs in
167 the ABruijn assembly showed links between chromosomes 3 and 13 (first contig),
168 chromosomes 3 and 2 (second contig), and chromosomes 10 and 2 (third contig). To verify
169 that the portions of these contigs were effectively chimeric, we back aligned the Nanopore
170 reads to the assembly and could not find any sequence that validated these links.
171 Unsurprisingly, these three chimeric contigs were fused at Ty1 transposable element
172 locations.

173 The alignment of each assembly to the reference genome showed that neither Canu, Miniasm,
174 nor SMARTdenovo could assemble the mitochondrial (Mt) genome completely. Because
175 ABruijn was the only assembler to assemble the complete Mt genome sequence, we decided
176 to use it to assemble the Mt DNA of the remaining 21 yeast strains (see below).

177 Generally, long reads allow tandem duplicated genes to be resolved, as for instance the *CUP1*
178 and *ENAI-2* gene families. We compared the maximum number of copies found in the
179 Nanopore reads and the estimated number of copies based on Illumina reads coverage of these
180 two tandem-repeated genes with the number of copies of these two genes in the four
181 assemblies (**Table S3**). After aligning the paired-end reads to the reference sequence and
182 computing of the coverage, we estimated that *CUP1* and *ENAI-2* were present in seven and
183 four copies, respectively. The maximum numbers of copies of these genes in a single
184 Nanopore read were eight for *CUP1* and five for *ENAI-2*. The numbers of copies of *CUP1*

185 and *ENAI-2* were, respectively, nine and three in the Canu assembly, seven and two in the
186 Miniasm assembly and seven and four in the SMARTdenovo and ABruijn assemblies.

187 The number of indels in each assembly was considerably high for each assembler. Thus, we
188 tested Nanopolish [2], the most commonly used Nanopore-only error corrector. We used the
189 SMARTdenovo assembly, which was the most continuous and gene-rich assembly and all 2D
190 reads for this test. After the error correction step, the cumulative length of the contigs
191 increased to 12.2 Mb and the N50 increased to 783 kb (at best it was 924 kb for the reference
192 genome). The number of mismatches, insertions and deletions decreased to 1,930, 7,707, and
193 17,445 respectively. The number of genes increased to 6,273 complete and 2,590 without an
194 indel (**Table 1**).

195 Although all metrics were improved, the number of indels was still too high, especially in the
196 coding regions of the genes. We decided to polish all assemblies with 2x250bp Illumina
197 paired-end reads at 300X genome coverage, using Pilon [28], to verify if the general quality
198 of the assembly improved. The polishing step increased the N50 of each assembly, and the
199 maximum of 816 kb was obtained with the ABruijn assembly. Pilon reduced the number of
200 errors of each assembly, and the Canu and ABruijn assemblies had the best base quality with
201 about 16 mismatches (15.85 and 17.88 for Canu and ABruijn respectively) and 22 indels
202 (22.49 and 21.76 for Canu and ABruijn respectively) per 100 kb. The SMARTdenovo
203 assembly contained the highest number of complete genes (6,266) and the Canu assembly
204 contained the highest number of genes without any indels (5,921) (**Table 2**). Furthermore, we
205 estimated the impact of the input coverage used to polish the consensus. We performed
206 successive polishing by using subsets of Illumina reads (ranging from 25X to 300X genome
207 coverage). We observed similar results in terms of number of mismatches and indels,
208 regardless of the input coverage. (**Figure S6**).

209 Finally, we evaluated the composition of each assembly for various elements (genes, repeated
210 elements, centromeres and telomeric regions). We also generated an Illumina-only assembly
211 using Spades assembler [29] to compare the number of features found in each assembly. All
212 the assemblies contained nearly the same number of centromeres (120 bp regions in the
213 reference genome assembly) and genes (**Figure 2**). The Nanopore assemblies contained more
214 complete genes than the Illumina one, however genes without indels are more frequent in the
215 Illumina-only assembly although nanopore assemblies were polished using Illumina reads.
216 even between 45 and 50 Long Terminal Repeat (LTR) retrotransposons (average size of 5.8
217 kb), while the Illumina-only assembly contained only one. The smallest number of telomeres
218 (three) was found in the ABruijn assembly, while nine, 18, 13, and 14 telomeres were found
219 in the Illumina, Canu, Miniasm, and SMARTdenovo assemblies, respectively. The Illumina-
220 only assembly contained five telomeric repeats (average size 100 bp), while the Nanopore-
221 only assemblies contained between six and nine telomeric repeats. The ABruijn assembly
222 contained the same number of genes encoded by the mitochondrial genome as the reference
223 sequence because it was the only assembler to fully assemble the Mt genome.

224 **S288C assemblies with R9 data**

225 The R9 version of the pore was released too late for us to use it to sequence all the natural *S.*
226 *cerevisiae* isolates. However, we did produce some data to compare the R7.3 and R9
227 assemblies. Because SMARTdenovo produced the best results (higher continuity and higher
228 gene content), we used it to assemble the R9 data generated from the S288C strain. We input
229 four different read datasets: all 1D and 2D reads, only 2D reads, 30x of the longest 2D reads
230 or 30x of the longest 1D and 2D reads (**Table S4**).

231 This time, the 30x of the longest 1D and 2D reads dataset gave the best results. Indeed, the
232 contiguity of the assembly increased, and the number of contigs decreased from 26 with the
233 R7.3 assembly to 23 with the R9 assembly. The number of indels also decreased from

234 133,676 with the R7.3 version to 95,012 with the R9 version. A direct consequence of using
1
2 235 the R9 version was that almost all the genes were found, and 6,302 of the 6,350 known genes
3
4
5 236 were complete and 1,226 did not contain any indels.
6
7

8 237 **Sequencing and assembly of the genomes of the 22 yeast strains**

9

10 238 To explore the variability of the genomic architecture within *S. cerevisiae*, 21 natural isolates
11
12 239 were sequenced in addition to the S288C reference genome using the same strategy, namely, a
13
14
15 240 combination of long Nanopore and short Illumina reads. Sequenced isolates were selected to
16
17 241 include as much diversity as possible in terms of global locations (including Europe, China,
18
19 242 Brazil, and Japan), ecological sources (such as fermented beverages, dairy products, trees and
20
21
22 243 fruit soil), as well as genetic variation highlighted in the frame of the extensive resequencing
23
24
25 244 1002 Yeast Genomes project (<http://1002genomes.u-strasbg.fr/>) (**Table S5**). Among these
26
27 245 isolates, the nucleotide variability was distributed across 491,076 segregating sites and the
28
29 246 genetic diversity, estimated by the average pairwise divergence (π), was 0.0062, which is
30
31
32 247 close to what is observed for the whole species [30].
33

34 248 A total of 78 MinION Mk1 runs were performed and the highest throughput we obtained was
35
36
37 249 650 Mb (1D and 2D reads). This led to 1.4 million of 2D reads with a cumulative length of 12
38
39 250 Gb. We obtained 2D coverage that ranged from 22x to 115x (**Figure S7**) among the strains
40
41
42 251 with a median read length of approximately 5.4 kb and a maximum size of 75 kb (**Figure S8**).
43
44 252 In general, three runs or less were sufficient to obtain the expected coverage. Next, for each
45
46
47 253 strain, we gave varying coverages of the longest 2D reads (**Table S6**) as input to
48
49 254 SMARTdenovo and retained the most contiguous assembly. These assemblies were then
50
51
52 255 given as input to Pilon for a polishing step with around 300x of Illumina paired-end reads
53
54 256 (each strain was individually sequenced using the Illumina technology). After polishing, we
55
56 257 obtained a median number of contigs of 27.5 (**Table 3**), the minimum number was for the CEI
57
58
59 258 strain (18 contigs) and the maximum was for the BAM strain (105 contigs). The median
60
61
62
63
64
65

259 cumulative length was 11.93 Mb and ranged from 11.83 Mb for the ADQ strain to 12.2 Mb
1
2 260 for the CNT strain. The median N50 contig size was 593 kb and varied from 201 kb for the
3
4 261 CIC strain to 896 kb for the ADQ strain. The L90 varied from 14 for the BCN, CEI, and CNT
5
6
7 262 strains, to 72 for the BAM strain with a median equal to 19.5.

8
9 263 To assemble the mitochondrial (Mt) genome, we used all the 2D reads as input to ABruijn. As
10
11
12 264 a result, we obtained an assembly for each strain and extracted the Mt genome after mapping
13
14 265 the contigs against the reference Mt genome. As was the case for the chromosomes, we used
15
16
17 266 Pilon with Illumina paired-end reads to obtain a corrected consensus sequence.

20 267 **Transposable elements**

21
22 268 The availability of high quality assemblies allowed us to establish an extensive map of the
23
24
25 269 transposable elements (TEs) to obtain a global view of their content and positions within the
26
27 270 21 natural yeast isolates (**Figure 3**). Using a reference sequence for each of the five known
28
29
30 271 TE families in yeast (namely Ty1 to Ty5), we mapped the TEs in each assembled genome.
31
32 272 Among the 50 annotated TEs in the S288C reference genome, 47 were detected at the correct
33
34
35 273 chromosomal locations in our assembly but three Ty1 locations were not recovered. Seven
36
37 274 additional Ty1 elements were found at unannotated sites, three of them have already been
38
39 275 detected in the reference genome [31]. These results attest to the high accuracy of our
40
41
42 276 assembly strategy for TE detection and localization. Among the 22 isolates, the TE content
43
44 277 was highly variable (**Table 4**), ranging from five to 55 elements, with a median value of 15.
45
46
47 278 While the frequency of the Ty4 and Ty5 elements was clearly low in all the isolates (up to
48
49 279 four and two elements, respectively), the Ty1, Ty2, and Ty3 elements were found in most of
50
51
52 280 the isolates. The most abundant TEs were Ty1 and Ty2, except in the Chinese BAM isolate,
53
54 281 in which 12 Ty3 elements were detected. As already described [32], the pattern of insertion of
55
56 282 these mobile elements is either specific to a given isolate, or shared by only a small number of
57
58
59 283 isolates (mostly two or three). However, four insertion hotspots have been highlighted (shared
60
61
62
63
64
65

284 by seven or more isolates) on chromosomes 2, 3, and 9. The shared insertion hotspots were
285 generally not specific to a specific Ty family, except for the hotspot located on a subtelomeric
286 region of the chromosome 3, which was specific to Ty5.

287 **Structural variations**

288 Structural variations such as copy number variants, large insertions and deletions,
289 duplications, inversions and translocations are of great importance at the phenotypic variation
290 level [33]. Compared with single nucleotide polymorphism (SNPs) and small indels, these
291 variants are usually more difficult to identify, in particular because resequencing strategies
292 have until recently focused mainly on the generation of short reads and reference-based
293 genome analysis. Nanopore long reads sequencing data allow the copy numbers of tandem
294 genes to be determined. As a testbed, we focused on two loci that are known to contain multi-
295 copy genes, namely *ENA* and *CUPI*. *ENA* genes encode plasma membrane Na⁺-ATPase
296 exporters, which play a role in the detoxification of Na⁺ ions in *S. cerevisiae*. *CUPI* genes
297 encode metallothioneins, which bind copper and are involved in resistance to copper exposure
298 by amplification of this locus. To determine the degree of divergence among the 21 strains,
299 we searched for the numbers of copies of the *CUPI* and *ENA*, two tandem-repeated genes in
300 the assemblies (**Table 5**). For this purpose, we extracted the corresponding sequence from the
301 S288C reference genome and aligned it to the assemblies of each strain. As expected and
302 already reported [34], the copy numbers of *ENA1* and *CUPI* varied greatly across the strains.
303 We found that the copy numbers of *ENA* genes in the 21 isolates ranged from 1 in 12 of the
304 genomes to five in the BHH strain (**Table 5**). The copy numbers of *CUPI* genes fluctuated
305 even more, ranging from one to 10 copies in the ABH and AEG strains. We also determined
306 the fitness of the 21 isolates in the presence of CuSO₄ and observed a correlation between the
307 number of *CUP* genes and the resistance of the strain to high concentration of CuSO₄ (**Figure**
308 **S9**).

309 Besides copy number variants, we also focused on larger structural variants, such as
1
2 310 translocations and inversions, because our highly contiguous assemblies allowed us to
3
4 311 investigate these events. We aligned the polished assemblies of the 21 strains to the reference
5
6
7 312 genome using NUCmer and inspected the alignments with the mummer software suite to
8
9
10 313 search for structural variations. We detected 29 translocations and four inversions within the
11
12 314 assemblies of 17 strains (**Table 6**). The median length of an inversion was 94 kb and their
13
14 315 breakpoints were located mostly in intergenic regions. It is well recognized that SVs might
15
16
17 316 play a major role in the genetic and phenotypic diversity in yeast [35, 36]. However, up to
18
19 317 now, it was impossible to assemble and have an exhaustive view of the SVs content in any *S.*
20
21
22 318 *cerevisiae* natural isolates. Indeed, short-read sequencing approaches are not suitable for SVs
23
24 319 studies because they results in a high number of false positive as well as false negative
25
26 320 detected events.

27
28
29 321 Among the detected events, one translocation detected between chromosomes 5 and 14 in the
30
31 322 ABH isolate and another translocation between chromosomes 7 and 12 in the AVB isolate
32
33
34 323 have already been described and confirmed in a reproductive isolation study in *S. cerevisiae*
35
36 324 [35]. A deeper investigation of our assemblies highlighted the presence of full-length Ty
37
38
39 325 transposons at some junctions of the translocation events. For example, the complex Ty-rich
40
41 326 junctions of the translocation between the chromosomes 7 and 12 in the ABH isolate was in
42
43
44 327 complete accordance with previously reported results [35]. Our results underline the high
45
46 328 resolution of the constructed assemblies, and show that complex events, such as
47
48
49 329 translocations, can be detected accurately with our strategy. Among the 22 isolates, six were
50
51 330 devoid of translocation events whereas the other 16 carries one to four such structural
52
53 331 rearrangements compared to the reference.

54
55
56 332 However, several limitations can be highlighted for these detections. Contrary to expectations,
57
58 333 no translocation that specifically affected subtelomeric regions was identified, underlining the
59
60
61
62
63
64
65

334 difficulty of discriminating regions that are variable and contain a large number of repeated
1
2 335 segments. Moreover, the detection accuracy is highly dependent on the completeness of the
3
4
5 336 assembly because, if translocation breakpoints are located on contigs boundaries, they will not
6
7 337 be detectable.
8
9

10 338 **Mitochondrial genome variation**

11
12
13 339 The ABruijn assembler allowed the construction of a single contig corresponding to the Mt
14
15 340 genome for each isolate. To assess the quality of the assemblies, we aligned the polished
16
17
18 341 S288C Mt contig to the reference sequence (GenBank: KP263414). Only four SNPs and few
19
20 342 indels, representing 15 bp of cumulative length, were detected. For all but two natural
21
22
23 343 isolates, all the Mt genes (eight protein coding genes, two rRNA subunits and 24 tRNAs)
24
25 344 were conserved and syntenous. The Mt genomes of the two remaining isolates (CNT and
26
27 345 CFF) contained one and two repeated regions covering a total of 6.5 and 8 kb, respectively. In
28
29
30 346 the CNT, the repeated region was in the *COXI* gene and affected its coding sequence. In the
31
32 347 CFF isolate, the *COXI*, *ATP6*, and *ATP8* genes would have been tandemly duplicated.
33
34
35 348 However, because we could not identify reads that clearly covered the repeated regions and
36
37 349 then confirmed the structural variations, we excluded these two Mt genome assemblies from
38
39
40 350 our dataset.
41

42 351 The sizes of the 20 considered assemblies ranged from 73.5 to 86.9 kb, which is close to the
43
44 352 size reported previously [37]. The differences in size between the assemblies can mainly be
45
46
47 353 attributed to the intron content of the *COXI* and *COB* genes (from two to eight introns in
48
49 354 *COXI* and from two to six introns in *COB*). These variations lead to extensive gene length
50
51
52 355 variability ranging from 5.7 kb to 14.9 kb for *COXI* and from 3.2kb to 8.6 kb for *COB*, while
53
54 356 the coding sequences of these 2 genes were exactly the same length among the 20 isolates.
55
56
57 357 Intergenic regions also accumulate many small indels, including those that affect the
58
59 358 interspersed GC-clusters, and a few large indels that sometimes correspond to variable
60
61
62
63
64
65

359 hypothetical open reading frames (ORFs), leading to sizes that range from 51.6 to 58 kb. To a
1
2 360 lesser extent, the 21S rRNA gene is also subjected to size variation that ranges from 3.2 to 4.4
3
4
5 361 kb.
6
7

8 362
9

12 363 Discussion

14 364 One of the major advantages of the Oxford Nanopore technology is the possibility of
15
16
17 365 sequencing very long DNA fragments. In our analyses, we obtained 2D reads up to 75 kb in
18
19 366 length, indicating that the system was able to read without interruption a flow of at least
20
21
22 367 150,000 nucleotides. Furthermore, the results of this analysis indicate that the error rate of the
23
24 368 ONT R7.3 reads was in the range that is obtained using existing long-read technologies (i.e.,
25
26 369 about 15% for 2D reads). However, the errors are not random and they significantly impact
27
28
29 370 stretches of the same nucleotides (homopolymers), which seems to be a feature inherent to the
30
31 371 ONT sequencing technology. Because the pore detects six nucleotides at a time, segmentation
32
33
34 372 of events is problematic in genomic regions with homopolymers longer than six bases [38].
35
36 373 With the current R7.3 release, homopolymers are prone to base deletion (representing 66% of
37
38
39 374 the errors observed in homopolymers). It may be improved with a steadier passing speed
40
41 375 through the pore or by increasing the speed of the molecule through the pore. In the same
42
43
44 376 way, the basecaller algorithm could be optimized to increase the accuracy per base. ONT have
45
46 377 recently reported several changes, including a fast mode (250 bp/second instead of 70
47
48
49 378 bp/second with R7.3 chemistry) and new basecaller software based on neural networks. These
50
51 379 new features are incorporated in the R9 version of MinION. We performed R9 experiments,
52
53 380 and observed a significant decrease in the error rate (with 1D and 2D reads, **Figure 1**). Using
54
55
56 381 this new release, homopolymers were more prone to base insertions (representing 63% of the
57
58 382 errors observed in homopolymers). Systematic errors are problematic for genome assembly
59
60
61
62
63
64
65

383 because they lead to the construction of less accurate consensus sequences. Furthermore,
1
2 384 indels negatively impact gene prediction because they can create frameshifts in the coding
3
4
5 385 regions of genes. We concluded that nanopore-only assemblies are difficult to use for analysis
6
7 386 at the gene level unless they are polished. However, polishing based only on nanopore reads
8
9
10 387 was not sufficient because although it reduced the number of indels by more than seven times,
11
12 388 we still had about 3,700 genes that were affected by potential frameshifts. The recently
13
14
15 389 developed R9 chemistry greatly improved the overall quality of the consensus sequences,
16
17 390 because starting with only 45x of 2D reads we obtained an assembly with the same contiguity
18
19 391 but with a decrease of nearly 30% in the number of indels (95,012 compared with 133,676).
20
21
22 392 We consider that the ONT sequencing platform will evolve in the coming years to produce
23
24 393 high quality long reads. Until then, a mixed strategy using high quality short reads remains the
25
26
27 394 only way to obtain high quality consensus sequences as well as a high level of contiguity.
28
29 395 Indeed, for the assembly of repetitive regions, the nanopore-only assemblies outperformed the
30
31 396 short-reads assemblies.

33
34 397 Our benchmark of nanopore-only assemblers shows that unfortunately a single “best
35
36 398 assembler” does not exist. Canu reconstructed the telomeric regions better and provided a
37
38
39 399 consensus of higher quality than Miniasm and SMARTdenovo. ABruijn seemed to produce
40
41 400 the most continuous assembly but some of the contigs were chimeric. However, ABruijn was
42
43
44 401 the only assembler to fully assemble the mitochondrial genome, and that is why we chose it to
45
46 402 assemble the Mt genomes of the 22 yeast strains. SMARTdenovo provided good overall
47
48
49 403 results for repetitive regions, completeness, contiguity, and speed. It was the most appropriate
50
51 404 choice to assemble the genome of all the yeast strains even if its major drawback was the
52
53 405 absence of the Mt genome sequence among the contig output.

55
56 406 The high contiguity of the 22 nanopore-only assemblies allowed us to detect transposable
57
58 407 element insertions and to provide a complete cartography of these elements. Ty1 was the most
59
60
61
62
63
64
65

408 abundant element and it was spread across the entire genome. Chromosome 12 was always
1
2 409 the most fragmented in our assemblies due to the presence of the rDNA cluster (around 100
3
4 410 copies in tandem). Furthermore, we easily identified known translocations (between
5
6
7 411 chromosomes 5 and 14 in the ABH isolate and between chromosomes 7 and 12 in the AVB
8
9 412 isolate). The high contiguity of the assemblies seemed to be limited by the read size rather
10
11
12 413 than the error rate. Work is still needed to prepare high-weight molecular DNA, enriched in
13
14 414 long fragments. The yeast genomes were successfully assembled with 8 kb and 20kb
15
16
17 415 fragment-sized libraries, but more complex genomes will require longer reads.
18
19

20 416 **Methods**

21 417 **DNA extraction**

22
23 418 Yeast cells were grown on YPD media (1% yeast extract, 2% peptone and 2% glucose) using
24
25 419 liquid culture or solid plates. Total genomic DNA was purified from 30 ml YPD culture using
26
27 420 Qiagen Genomic-Tips 100/G and Genomic DNA Buffers as per the manufacturer's
28
29 421 instructions. The quantity and quality of the extracted DNA were controlled by migration on
30
31
32 422 agarose gel, spectrophotometry (NanoDrop ND-1000, ThermoFisher, Wilmington, DE, USA),
33
34 423 and fluorometric quantification (Qubit, ThermoFisher, Wilmington, USA).
35
36
37
38

39 424 **Illumina PCR-free library preparation and sequencing**

40
41
42 425 DNA (6 µg) was sonicated to a 100 to 1500 bp size range using a Covaris E210 sonicator
43
44 426 (Covaris, Woburn, MA, USA). Fragments were end-repaired using the NEBNext® End
45
46 427 Repair Module (New England Biolabs, Ipswich, MA, USA) and 3'-adenylated with the
47
48
49 428 NEBNext dA-Tailing Module. Illumina adapters were added using the NEBNext Quick
50
51
52 429 Ligation Module. Ligation products were purified with AMPure XP beads (Beckmann
53
54 430 Coulter Genomics, Danvers, MA, USA). Libraries were quantified by qPCR using the KAPA
55
56 431 Library Quantification Kit for Illumina Libraries (KapaBiosystems, Wilmington, MA, USA)
57
58
59 432 and library profiles were assessed using a DNA High Sensitivity LabChip kit on an Agilent
60
61
62
63
64
65

433 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). Libraries were sequenced on an
1
2 434 Illumina MiSeq or a HiSeq 2500 instrument (San Diego, CA, USA) using 300 or 250 base-
3
4
5 435 length read chemistry in a paired-end mode.

7 436 **Nanopore 20 kb libraries preparation**

8
9 437 MinION sequencing libraries were prepared according to the SQK-MAP005 or SQK-
10
11 438 MAP006-MinION gDNA Sequencing Kit protocols. Six to 10 µg of genomic DNA was
12
13
14 439 sheared to approximately 20,000 bp with g-TUBE (Covaris, Woburn, MA, USA). After clean-
15
16
17 440 up using 0.4x AMPure XP beads, sequencing libraries were prepared according to the SQK-
18
19 441 MAP005 or SQK-MAP006 Sequencing Kit protocols, including the PreCR treatment (NEB,
20
21
22 442 Ipswich, USA) for the SQK-MAP005 protocol or the NEBNext FFPE DNA repair step (NEB,
23
24 443 Ipswich, USA) for the SQK-MAP006 protocol.

26 444 **Nanopore 8 kb libraries preparation**

27
28
29 445 MinION sequencing libraries were prepared according to the SQK-MAP005 or SQK-
30
31 446 MAP006-MinION gDNA Sequencing Kit protocols. Two µg of genomic DNA was sheared to
32
33
34 447 approximately 8,000 bp with g-TUBE. After clean-up using 1x AMPure XP beads,
35
36 448 sequencing libraries were prepared according to the SQK-MAP005 or SQK-MAP006
37
38
39 449 Sequencing Kit protocol, including the PreCR treatment for the SQK-MAP005 protocol or the
40
41 450 NEBNext FFPE DNA repair step for the SQK-MAP006 protocol.

43 451 **Nanopore Low input 8 kb libraries preparation**

44
45
46 452 The following protocol was applied to some samples (**Supplementary File 3**). Five hundred
47
48
49 453 ng of genomic DNA was sheared to approximately 8,000 bp with g-TUBE. After clean-up
50
51 454 using 1x AMPure XP beads and the NEBNext FFPE DNA repair step, 100 ng of DNA was
52
53 455 prepared according to the Low Input Expansion Pack Protocol for genomic DNA.

56 456 **MinION™ flow cell preparation and sample loading**

1 457 The sequencing mix was prepared with 8 μ L of the DNA library, water, the Fuel Mix and the
2 458 Running buffer according to the SQK-MAP005 or the SQK-MAP006 protocols. The
3
4 459 sequencing mix was added to the R7.3 flowcell for a 48 hours run. The flowcell was then
5
6
7 460 reloaded three times according to the following schedule: 5 hours (4 μ L of DNA library), 24
8
9 461 hours (8 μ L of DNA library) and 29 hours (4 μ L of DNA library). Regarding the Low Input
10
11 462 libraries, the flowcell was loaded and then reloaded after 24 hours of run time with a
12
13 463 sequencing mix containing 10 μ L of the DNA library (**Supplementary File 3**).

16 464 **MinION[®] sequencing and reads filtering**

17
18
19 465 Read event data generated by MinKNOW[™] control software (version 0.50.1.15 to 0.51.1.62)
20
21 466 were base-called using the Metrichor[™] software (version 2.26.1 to 2.38.3). The data
22
23 467 generated (pores metrics, sequencing, and base-calling data) by MinION software were stored
24
25 468 and organized using a Hierarchical Data Format (HDF5). Three types of reads were obtained:
26
27 469 template, complement, and two-directions (2D). The template and complement reads
28
29 470 correspond to sequencing of the two DNA strands. Metrichor combines template and
30
31 471 complement reads to produce a consensus (2D) sequence [39]. FASTA reads were extracted
32
33 472 from MinION HDF5 files using poretools [40]. To assess the quality of the MinION reads, we
34
35 473 aligned reads against the *S. cerevisiae* S288C reference genome using the LAST aligner
36
37 474 (version 588) [41]. Because the MinION reads are long and have a high error rate we used a
38
39 475 gap open penalty of 1 and a gap extension penalty of 1.

46 476 **Illumina reads processing and quality filtering**

47
48 477 After the Illumina sequencing, an in-house quality control process was applied to the reads
49
50 478 that passed the Illumina quality filters. The first step discards low-quality nucleotides ($Q < 20$)
51
52 479 from both ends of the reads. Next, Illumina sequencing adapters and primers sequences were
53
54 480 removed from the reads. Then, reads shorter than 30 nucleotides after trimming were
55
56 481 discarded. These trimming and removal steps were achieved using in-house-designed
57
58
59
60
61
62
63
64
65

1 482 software based on the FastX package [42]. The last step identifies and discards read pairs that
2 483 mapped to the phage phiX genome, using SOAP [43] and the phiX reference sequence
3
4 484 (GenBank: NC_001422.1). This processing resulted in high-quality data and improvement of
5
6
7 485 the subsequent analyses.
8

9 486 **Assembler evaluation**

10
11
12 487 To determine the assembler to use on the *de novo* sequenced 22 yeast strains, tests were
13
14 488 conducted on S288C, the only *S. cerevisiae* strain for which there is an established reference
15
16
17 489 genome. We used different subsets of the reads as input to Canu (github commit ae9eccc),
18
19 490 Miniasm (github commit 17d5bd1), SMARTdenovo (github commit 61cf13d), and ABruijn
20
21 491 (github commit dc209ee), four assemblers that can take advantage of long reads. These
22
23
24 492 subsets consisted of varying coverages of 1D, 2D, 2D pass reads, which are 2D reads that
25
26
27 493 have an average quality greater than nine, and reads corrected by Canu. Canu was executed
28
29 494 with the following parameters: genomeSize=12m, minReadLength=5000,
30
31 495 mhapSensitivity=high, corMhapSensitivity=high, errorRate=0.01 and corOutCoverage=500.
32
33
34 496 Miniasm was run with the default parameters indicated on the github website. SMARTdenovo
35
36 497 was executed with the default parameters and -c 1 to run the consensus step. ABruijn was run
37
38
39 498 with default parameters. After the assembly step, we polished each set of contigs with Pilon
40
41 499 (version 1.1.12) using 300X of Illumina 2x250 bp paired-end reads. Assemblies were aligned
42
43
44 500 to the S288C reference genome using Quast in conjunction with the GFF file of S288C to
45
46 501 detect assembly errors, and complete and partial genes. We also visualized the alignments
47
48
49 502 using mummerplot to detect chimeric contigs.
50

51 503 **Assembly of the genome of the 22 yeast strains**

52
53 504 The 22 genomes were assembled by utilizing varying sequencing coverage, going from 10X
54
55
56 505 to 50X, of the longest 2D reads as input to SMARTdenovo with the default parameters and -c
57
58 506 1 to run the consensus step. Then, for each strain, the most contiguous assembly (based on the
59
60
61
62
63
64
65

507 N50 and the number of contigs) was polished using ~300X of 2x250bp Illumina paired-end
1
2 508 reads (each yeast strain was sequenced separately beforehand).
3

4 509 **Genes and transposons detection**

5
6
7 510 To detect genes and transposons in the assemblies, we extracted the corresponding sequences
8
9 511 from the reference genome. We then mapped these elements to the assemblies using the Last
10
11 512 aligner. Only alignments that showed more than 80% identity over at least 90% of the
12
13 513 sequence length were retained and considered as a match. We used a similar procedure to
14
15 514 count the maximum number of gene in the Nanopore reads dataset, the only modification was
16
17 515 that the percentage identity had to be at least 70% to account for the high error rate of the
18
19 516 reads. To estimate the number of copies in the Illumina reads, we aligned paired-end reads to
20
21 517 the reference genome with BWA aln and then computed the coverage using samtools mpileup
22
23 518 algorithm [44] and divided the number we obtained for each region of interest by the median
24
25 519 coverage of the corresponding chromosome.
26
27
28
29
30

31 520 **Feature number estimation**

32
33
34 521 We generated an Illumina-only assembly using Spades version v3.7.0 with default parameters
35
36 522 and compare the completeness of this assembly to the nanopore-only assemblies. To estimate
37
38 523 the number of features across all S288C assemblies, we aligned each post-polishing consensus
39
40 524 sequence to the S288C reference genome using NUCmer. Only the best alignments were
41
42 525 conserved by using the *delta-filter -l* command. Next, we used the bedtools suite [45] with
43
44 526 the command *bedtools intersect -u -wa -f 0.99* to compare the alignments to the reference
45
46 527 GFF file. Finally, we counted the number of features of our interest.
47
48
49
50

51 528 **Circularization of mitochondrial genomes**

52
53 529 To circularize the Mt genomes, we split the contig corresponding to the Mt sequence in each
54
55 530 strain into two distinct contigs. Then, we gave the two contigs as input to the minimus2 [46]
56
57 531 tool from the AMOS package. As a result, we obtained a single contig that did not contain the
58
59
60
61
62
63
64
65

532 overlap corresponding to the circularization zone. Finally, to start the Mt sequence of all
1
2 533 isolates at the same position as the reference, we mapped each Mt sequence to the reference
3
4 534 using NUCmer. The *show-coords* command allowed us to identify the position in the Mt
5
6
7 535 sequences of all the strains that corresponded to the first position of the reference Mt genome.
8
9 536

537 **Declarations**

538 **Availability of Data and Materials**

539 The 22 genome assemblies are freely available at <http://www.genoscope.cns.fr/yeast>. The
17
18 540 Illumina and MinION data are available in the European Nucleotide Archive under accession
19
20
21 541 number ERP016443. Supporting data is also available from the *GigaScience* GigaDB
22
23 542 repository [47].
24
25

544 **Abbreviations**

545 ONT: Oxford Nanopore Technology; SMRT: Single-Molecule Real-Time Sequencing; USB:
31
32 546 Universal Serial Bus; Mt: Mitochondrial; LTR: Long Terminal Repeat; SNP: Single
33
34
35 547 Nucleotide Polymorphism; ORF: Open Reading Frame; MAP: MinION Access Programme.
36
37

549 **Ethics approval and consent to participate**

550 Not applicable
41
42
43

552 **Consent for publication**

553 Not applicable
47
48
49

555 **Competing interests**

556 The authors declare that they have no competing interests. Oxford Nanopore Technologies
54
55 557 Ltd contributed to this study by providing some of the R9 reagents free of charges. BI, SD,
56
57
58 558 CCR, AL, SE, PW and JMA are part of the MinION Access Programme (MAP).
59
60

560 **Funding**

1 561 This work was supported by the Genoscope, the Commissariat à l’Energie Atomique et aux
2
3 562 Energies Alternatives (CEA), France Génomique (ANR-10-INBS-09-08) and the Agence
4
5 563 Nationale de la Recherche (ANR-16-CE12-0019).
6

7 564
8

9 **Author’s contributions**

10 565 CCA extracted the DNA. EP, OB, CCR and AL optimized and performed the sequencing. BI,
11
12
13 567 AF, LDA, SF, SD, SE and JMA performed the bioinformatic analyses. BI, AF, JS and JMA
14
15
16 568 wrote the article. GL, PW, JS and JMA supervised the study.
17
18

19 569
20

21 **Acknowledgements**

22
23 571 The authors are grateful to Oxford Nanopore Technologies Ltd for providing early access to
24
25 572 the MinION device through the MinION Access Programme (MAP) and we thank the staff of
26
27
28 573 Oxford Nanopore Technology Ltd for technical help. The authors acknowledge Pierre Le Ber
29
30 574 and Claude Scarpelli for continuous support. JS is a member of the Institut Universitaire de
31
32
33 575 France.
34

35 576
36

37 **Additional files**

38 577
39 578 All the supporting data are included as a three additional files: a first one which contains
40
41
42 579 Figures S1-S9 and Tables S1-S6 and two excel files which contain the metrics of all
43
44 580 assemblies generated in this study and the description of each MinION run.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

References

1. Jain M, Fiddes IT, Miga KH, Olsen HE, Paten B, Akeson M: **Improved data analysis for the MinION nanopore sequencer.** *Nature methods* 2015, **12**(4):351-356.
2. Loman NJ, Quick J, Simpson JT: **A complete bacterial genome assembled de novo using only nanopore sequencing data.** *Nature methods* 2015, **12**(8):733-735.
3. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE *et al*: **Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data.** *Nature methods* 2013, **10**(6):563-569.
4. Koren S, Phillippy AM: **One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly.** *Current opinion in microbiology* 2015, **23**:110-120.
5. Koren S, Schatz MC, Walenz BP, Martin J, Howard JT, Ganapathy G, Wang Z, Rasko DA, McCombie WR, Jarvis ED *et al*: **Hybrid error correction and de novo assembly of single-molecule sequencing reads.** *Nature biotechnology* 2012, **30**(7):693-700.
6. Madoui MA, Engelen S, Cruaud C, Belser C, Bertrand L, Alberti A, Lemainque A, Wincker P, Aury JM: **Genome assembly using Nanopore-guided long and error-free DNA reads.** *BMC genomics* 2015, **16**:327.
7. Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz MC, McCombie WR: **Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a eukaryotic genome.** *Genome research* 2015, **25**(11):1750-1756.
8. Huddleston J, Ranade S, Malig M, Antonacci F, Chaisson M, Hon L, Sudmant PH, Graves TA, Alkan C, Dennis MY *et al*: **Reconstructing complex regions of genomes using long-read sequencing technology.** *Genome research* 2014, **24**(4):688-696.
9. Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, Antonacci F, Surti U, Sandstrom R, Boitano M *et al*: **Resolving the complexity of the human genome using single-molecule sequencing.** *Nature* 2015, **517**(7536):608-611.
10. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM *et al*: **Haplotyping germline and cancer genomes with high-throughput linked-read sequencing.** *Nature biotechnology* 2016, **34**(3):303-311.
11. Mostovoy Y, Levy-Sakin M, Lam J, Lam ET, Hastie AR, Marks P, Lee J, Chu C, Lin C, Dzakula Z *et al*: **A hybrid approach for de novo human genome sequence assembly and phasing.** *Nature methods* 2016, **13**(7):587-590.
12. Weisenfeld NI, Kumar V, Shah P, Church D, Jaffe DB: **Direct determination of diploid genome sequences.** *bioRxiv* 2016.
13. Mardis ER: **Next-generation DNA sequencing methods.** *Annual review of genomics and human genetics* 2008, **9**:387-402.
14. Loman NJ, Watson M: **Successful test launch for nanopore sequencing.** *Nature methods* 2015, **12**(4):303-304.
15. Deamer D, Akeson M, Branton D: **Three decades of nanopore sequencing.** *Nature biotechnology* 2016, **34**(5):518-524.
16. Kasianowicz JJ, Brandin E, Branton D, Deamer DW: **Characterization of individual polynucleotide molecules using a membrane channel.** *Proceedings of the National Academy of Sciences of the United States of America* 1996, **93**(24):13770-13773.

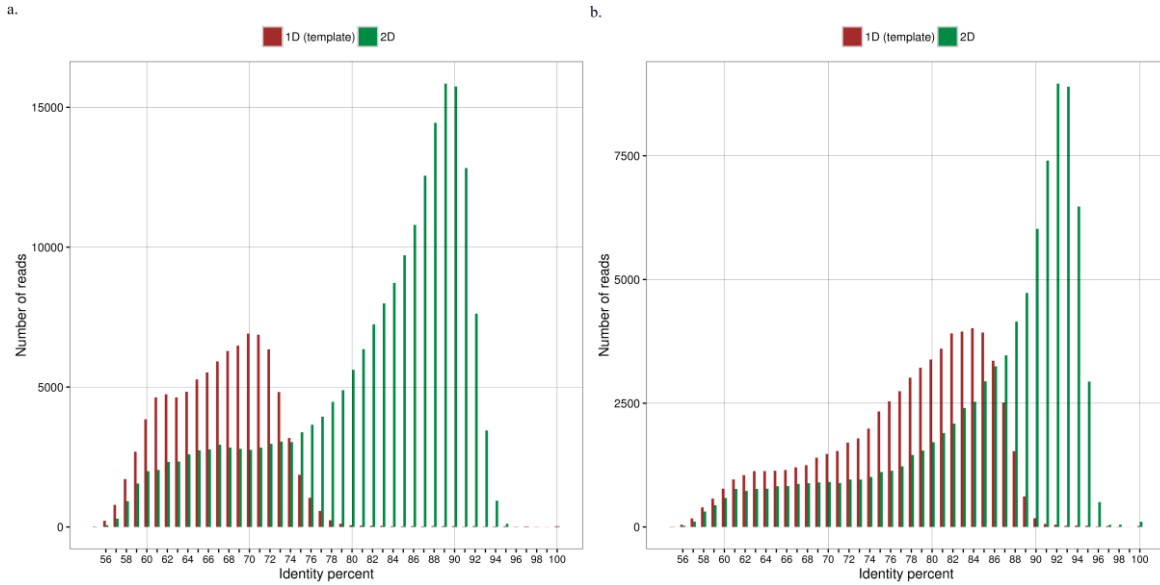
17. Cherf GM, Lieberman KR, Rashid H, Lam CE, Karplus K, Akeson M: **Automated forward and reverse ratcheting of DNA in a nanopore at 5-A precision.** *Nature biotechnology* 2012, **30**(4):344-348.
18. Manrao EA, Derrington IM, Laszlo AH, Langford KW, Hopper MK, Gillgren N, Pavlenok M, Niederweis M, Gundlach JH: **Reading DNA at single-nucleotide resolution with a mutant MspA nanopore and phi29 DNA polymerase.** *Nature biotechnology* 2012, **30**(4):349-353.
19. Laszlo AH, Derrington IM, Ross BC, Brinkerhoff H, Adey A, Nova IC, Craig JM, Langford KW, Samson JM, Daza R *et al*: **Decoding long nanopore sequencing reads of natural DNA.** *Nature biotechnology* 2014, **32**(8):829-833.
20. Loose M, Malla S, Stout M: **Real-time selective sequencing using nanopore technology.** *Nature methods* 2016.
21. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform.** *Bioinformatics* 2009, **25**(14):1754-1760.
22. Berlin K, Koren S, Chin CS, Drake JP, Landolin JM, Phillippy AM: **Assembling large genomes with single-molecule sequencing and locality-sensitive hashing.** *Nature biotechnology* 2015, **33**(6):623-630.
23. Li H: **Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences.** *Bioinformatics* 2016, **32**(14):2103-2110.
24. <https://github.com/ruanjue/smartdenovo>
25. Lin Y, Yuan J, Kolmogorov M, Shen MW, Pevzner PA: **Assembly of Long Error-Prone Reads Using de Bruijn Graphs.** *bioRxiv* 2016.
26. Gurevich A, Saveliev V, Vyahhi N, Tesler G: **QUAST: quality assessment tool for genome assemblies.** *Bioinformatics* 2013, **29**(8):1072-1075.
27. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL: **Versatile and open software for comparing large genomes.** *Genome biology* 2004, **5**(2):R12.
28. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J, Young SK *et al*: **Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement.** *PloS one* 2014, **9**(11):e112963.
29. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Prjibelski AD *et al*: **SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing.** *Journal of computational biology : a journal of computational molecular cell biology* 2012, **19**(5):455-477.
30. Peter J, Schacherer J: **Population genomics of yeasts: towards a comprehensive view across a broad evolutionary scale.** *Yeast* 2016, **33**(3):73-81.
31. Bleykasten-Grosshans C, Jung PP, Fritsch ES, Potier S, de Montigny J, Souciet JL: **The Ty1 LTR-retrotransposon population in *Saccharomyces cerevisiae* genome: dynamics and sequence variations during mobility.** *FEMS yeast research* 2011, **11**(4):334-344.
32. Bleykasten-Grosshans C, Friedrich A, Schacherer J: **Genome-wide analysis of intraspecific transposon diversity in yeast.** *BMC genomics* 2013, **14**:399.
33. Weischenfeldt J, Symmons O, Spitz F, Korbel JO: **Phenotypic impact of genomic structural variation: insights from and for human disease.** *Nature reviews Genetics* 2013, **14**(2):125-138.
34. Strobe PK, Skelly DA, Kozmin SG, Mahadevan G, Stone EA, Magwene PM, Dietrich FS, McCusker JH: **The 100-genomes strains, an *S. cerevisiae* resource that illuminates its natural phenotypic and genotypic variation and emergence as an opportunistic pathogen.** *Genome research* 2015, **25**(5):762-774.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
35. Hou J, Friedrich A, de Montigny J, Schacherer J: **Chromosomal rearrangements as a major mechanism in the onset of reproductive isolation in *Saccharomyces cerevisiae***. *Current biology : CB* 2014, **24**(10):1153-1159.
 36. Naseeb S, Carter Z, Minnis D, Donaldson I, Zeef L, Delneri D: **Widespread Impact of Chromosomal Inversions on Gene Expression Uncovers Robustness via Phenotypic Buffering**. *Molecular biology and evolution* 2016, **33**(7):1679-1696.
 37. Wolters JF, Chiu K, Fiumera HL: **Population structure of mitochondrial genomes in *Saccharomyces cerevisiae***. *BMC genomics* 2015, **16**:451.
 38. David M, Dursi LJ, Yao D, Boutros PC, Simpson JT: **Nanocall: An Open Source Basecaller for Oxford Nanopore Sequencing Data**. *bioRxiv* 2016.
 39. Quick J, Quinlan AR, Loman NJ: **A reference bacterial genome dataset generated on the MinION portable single-molecule nanopore sequencer**. *GigaScience* 2014, **3**:22.
 40. Loman NJ, Quinlan AR: **Poretools: a toolkit for analyzing nanopore sequence data**. *Bioinformatics* 2014, **30**(23):3399-3401.
 41. Kielbasa SM, Wan R, Sato K, Horton P, Frith MC: **Adaptive seeds tame genomic sequence comparison**. *Genome research* 2011, **21**(3):487-493.
 42. http://hannonlab.cshl.edu/fastx_toolkit/
 43. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J: **SOAP2: an improved ultrafast tool for short read alignment**. *Bioinformatics* 2009, **25**(15):1966-1967.
 44. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R: **The Sequence Alignment/Map format and SAMtools**. *Bioinformatics* 2009, **25**(16):2078-2079.
 45. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic features**. *Bioinformatics* 2010, **26**(6):841-842.
 46. Schatz MC, Phillippy AM, Sommer DD, Delcher AL, Puiu D, Narzisi G, Salzberg SL, Pop M: **Hawkeye and AMOS: visualizing and assessing the quality of genome assemblies**. *Briefings in bioinformatics* 2013, **14**(2):213-224.
 47. Istace, B; Friedrich, A; d'Agata, L; Faye, S; Payen, E; Beluche, O; Caradec, C; Davidas, S; Cruaud, C; Liti, G; Lemainque, A; Engelen, S; Wincker, P; Schacherer, J; Aury, J, M (2016): **Supporting data for "de novo assembly and population genomic survey of natural yeast isolates with the Oxford Nanopore MinION sequencer"** *GigaScience Database*. <http://dx.doi.org/10.5524/100263>

581 **Figures**

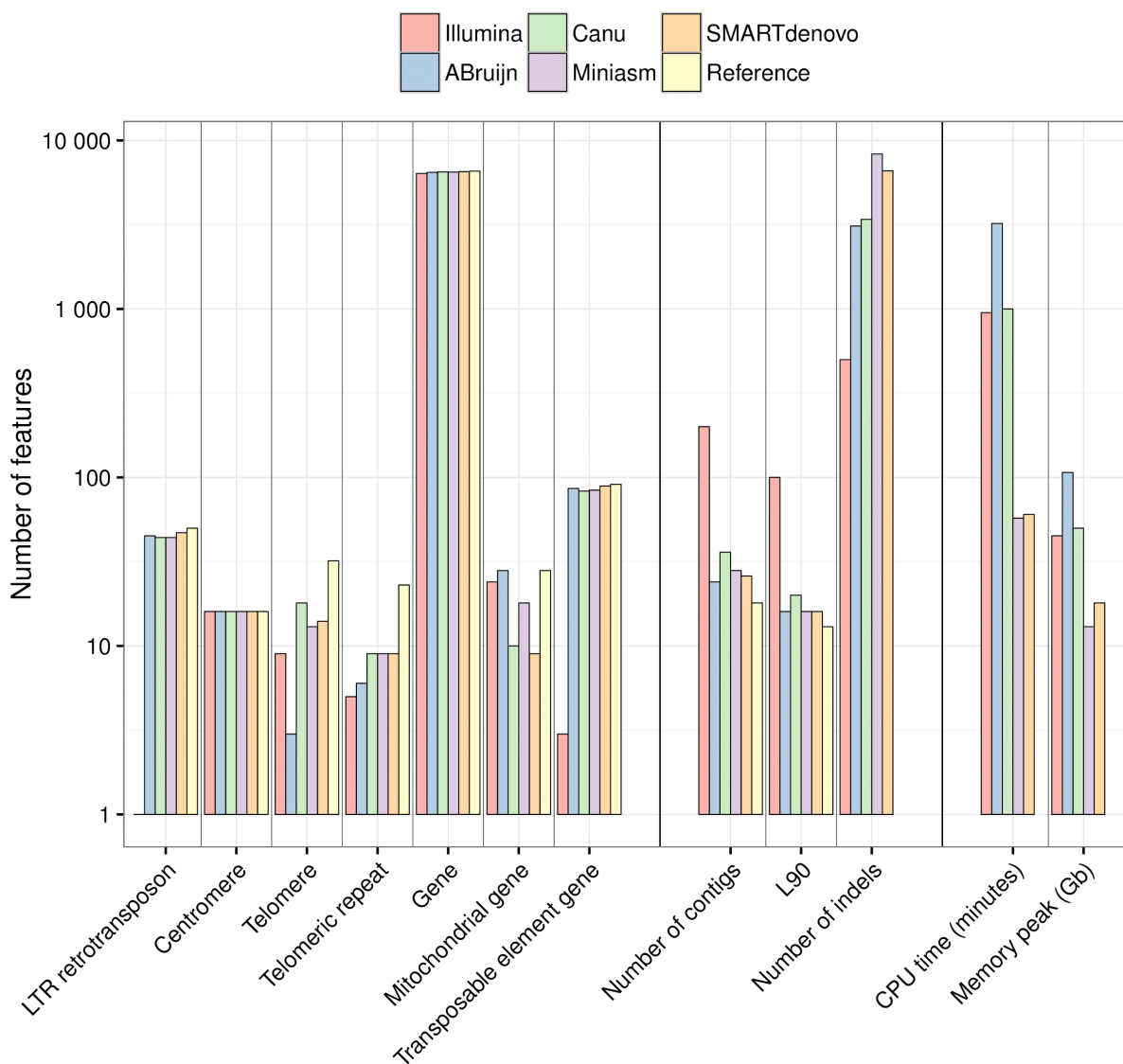
582 **Figure 1: Identity distribution of Nanopore reads.** Percent identity of the aligned MinION
583 1D (red bars) and 2D (green bars) reads. The MinION reads were aligned using LAST
584 software. **a.** R7.3 chemistry **b.** R9 chemistry

585



586

587 **Figure 2: Feature composition of the S288C assemblies, assembly and quality metrics**
 588 **and assembler running statistics.** The feature content of the best S288C assemblies for each
 589 assembler is shown in the left part of the figure. The feature composition was obtained by
 590 aligning each assembly to the S288C reference genome. Assembly and quality metrics for
 591 each assembly, obtained by using Quast, are shown in the middle part of the figure. The
 592 running time and the memory usage of each assembler are shown in the right part of the
 593 figure.



594

595 **Figure 3: Cartography of the Ty transposon family.** First and second tracks show,
 1
 2 596 respectively, the percentage identity of the SMARTdenovo S288C assembly before and after
 3
 4 597 polishing with Illumina paired-end reads using Pilon. The third track shows the 80th percentile
 5
 6 598 number of contigs obtained for each strain and for all chromosomes. The remaining tracks
 7 599 show the density of Ty transposons or positions of the Ty1, Ty2, Ty3, Ty4, and Ty5
 8
 9 600 transposons across all the yeast strains. The red dot on the karyotype track shows the position
 10
 11 601 of the rDNA cluster.



602

603 **Table 1: Metrics of the SMARTdenovo S288C assemblies before and after polishing**
 1 604 **with Nanopolish using R7 reads.** The Nanopore 2D reads were aligned to the most
 2 605 continuous SMARTdenovo assembly. The alignment was given as input to Nanopolish to
 3 606 correct assembly errors. Metrics were obtained by aligning the pre-polishing and post-
 4 607 polishing version of the assembly to the reference genome using Quast.

	SMARTdenovo Pre-polishing	SMARTdenovo Post-polishing
# contigs	26	26
Cumulative size	12,018,244	12,204,373
N50	771,149	782,423
N90	238,808	242,444
L50	7	7
L90	16	16
# mismatches	6,970	1,930
# insertions	7,735	7,707
# deletions	128,050	17,445
# deletions in homopolymers	79,152	6,869
# genes	6,251 + 24 partial	6,273 + 15 partial
# genes without indels	429	2,590

608

609 **Table 2: Metrics of the S288C assemblies after polishing.** Assemblies were corrected using
 1 610 300x of 2x250bp Illumina reads as input to Pilon. The resulting corrected assembly was then
 2 611 aligned to the S288C reference genome using Quast.

	Spades	Canu	Miniasm	SMARTdenovo	ABruijn
Reads dataset used	Illumina PE 2x250 bp	2D pass	Canu-corrected	Longest 2D	2D
Coverage	300x	67x	108x	30x	120x
# reads > 10kb	0	16,860	21,005	28,668	28,668
# contigs	376	37	28	26	23
Cumulative size	12,047,788	12,230,747	12,113,521	12,213,590	12,182,847
Genome fraction (%)	96.464	98.519	98.421	99.352	98.635
N50	149,184	610,494	736,456	783,336	816,355
N90	19,522	191,846	265,917	242,658	257,117
L50	27	8	7	7	6
L90	100	20	16	16	16
# mismatches	1,126	1,898	4,455	4,205	2,138
# mismatches per 100 kb	9.47	15.85	37.23	34.27	17.88
# insertions	81	1,657	3,164	2,384	1,325
# deletions	439	1,869	5,208	5,551	1,838
# deletions in homopolymers	38	868	4,248	4,023	740
#indels per 100 kb	1.97	22.49	57.27	46.76	21.76
# genes	6,087 + 177 partial	6,241 + 32 partial	6,215 + 37 partial	6,266 + 33 partial	6,243 + 45 partial
# genes without indels	6,023	5,921	5,475	5,881	6,002

612

613 **Table 3: Assembly metrics of the SMARTdenovo assemblies of all yeast strain genomes.**

	# contigs	Cumul (bp)	N50 (bp)	N90 (bp)	L50	L90	Max size (bp)
ABH	22	11,960,929	803,880	267,734	6	16	1,483,918
ADM	41	11,883,044	474,542	171,488	10	26	1,009,064
ADQ	26	11,828,347	896,166	223,992	6	18	1,223,692
ADS	33	11,706,636	524,733	247,699	9	21	1,050,223
AEG	23	12,026,175	681,360	273,814	7	16	1,244,014
AKR	25	11,911,766	729,090	243,900	7	17	1,056,085
ANE	47	11,900,397	312,705	144,286	11	31	933,716
ASN	40	11,904,493	394,798	143,405	11	28	846,371
AVB	31	11,991,127	609,633	199,011	7	20	1,225,549
BAH	28	11,829,394	571,862	227,561	8	20	1,066,359
BAL	27	11,907,375	678,155	269,114	7	19	1,075,839
BAM	105	11,996,380	162,412	53,623	24	72	450,388
BCN	19	11,775,292	785,507	458,793	6	14	1,410,650
BDF	45	12,068,568	460,458	116,953	10	29	863,099
BHH	26	11,973,506	577,727	221,661	7	18	1,530,377
CBM	68	11,553,446	258,798	86,167	16	44	521,412
CEI	18	11,987,201	800,227	451,575	6	14	1,480,681
CFA	24	11,834,226	726,317	225,716	7	17	1,032,352
CFF	81	12,162,869	236,957	83,285	18	54	550,022
CIC	96	12,016,445	201,870	63,799	22	63	377,026
CNT	22	12,171,929	800,046	440,742	6	14	1,402,970
CRV (S288C)	26	12,213,584	783,337	242,658	7	16	1,532,642
Median	27.5	11,936,347	593,680	224,854	7	19.5	1,061,222
Reference	17	12,157,105	924,431	439,888	6	13	1,531,933

614

615 **Table 4: Number of copies of multiple transposons across all yeast strains assemblies.**

	Ty1	Ty2	Ty3	Ty4	Ty5
ABH	4	7	6	3	2
ADM	5	8	1	1	0
ADQ	4	7	1	2	0
ADS	1	9	0	0	1
AEG	15	7	2	1	2
AKR	4	4	4	1	1
ANE	1	5	3	2	0
ASN	13	6	0	0	0
AVB	0	29	0	0	2
BAH	0	6	1	3	0
BAL	8	0	12	0	0
BAM	4	13	6	2	1
BCN	6	0	0	0	0
BDF	13	3	3	3	1
BHH	20	12	5	4	0
CBM	3	1	0	1	0
CEI	2	20	1	0	0
CFA	8	1	1	0	1
CFF	6	6	2	0	1
CIC	6	3	1	1	0
CNT	17	6	1	1	1
CRV (S288C)	36	13	2	3	1
Reference	31	13	2	3	1

616

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

617 **Table 5: Copy number of *CUP1* and *ENA1-2* tandem-repeated genes across the 21**
 618 **natural isolates assemblies.**

	<i>ENA1-2</i>	<i>CUP1</i>
ABH	1	10
ADM	2	1
ADQ	1	1
ADS	2	3
AEG	2	10
AKR	1	1
ANE	1	1
ASN	1	3
AVB	4	2
BAH	1	1
BAL	1	1
BAM	1	2
BCN	1	1
BDF	4	4
BHH	5	3
CBM	1	1
CEI	1	1
CFA	1	1
CFF	2	4
CIC	2	4
CNT	2	1

619

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

620 **Table 6: Chromosomic rearrangements detected across all 21 strains.**

Strain	Chromosome 1	Chromosome 2	Type
ABH	5	14	Translocation
ABH	5	14	Translocation
ABH	5	14	Translocation
ABH	14	14	Inversion
ADM	2	4	Translocation
ADM	5	7	Translocation
AKR	15	4	Translocation
ANE	16	5	Translocation
ANE	9	14	Translocation
ASN	5	2	Translocation
AVB	12	7	Translocation
AVB	7	12	Translocation
BAH	4	7	Translocation
BAH	10	9	Translocation
BAL	8	9	Translocation
BAM	4	7	Translocation
BAM	12	13	Translocation
BCN	6	13	Translocation
BCN	6	15	Translocation
BDF	4	14	Translocation
BDF	4	4	Inversion
BDF	5	12	Translocation
BDF	10	5	Translocation
BHH	12	12	Inversion
BHH	12	12	Inversion
CBM	16	3	Translocation
CBM	4	7	Translocation
CBM	12	15	Translocation
CEI	11	12	Translocation
CFF	14	12	Translocation
CIC	11	8	Translocation
CIC	4	7	Translocation
CNT	6	14	Translocation

622



Click here to access/download
Supplementary Material
Supplementary_File1.docx





Click here to access/download
Supplementary Material
Supplementary_File2.xlsx



Click here to access/download
Supplementary Material
Supplementary_File3.xlsx

Reviewer's report

Title: de novo assembly and population genomic survey of natural yeast isolates with the Oxford Nanopore MinION sequencer

Reviewer number: 1

Reviewer's report:

The authors describe assembly experiments on a set of yeast isolates using Oxford Nanopore MinION technology, both the older (and now discontinued) R7.3 and the newer (but about to be superseded) R9 chemistries. The methods are well-described, the data has been deposited in a stable archive and the results section performs a number of useful assessments of the quality of the assemblies using different de novo assembly tools.

Table S2 appears to be a subset of the columns of Table 1. If it doesn't provide any additional information, it should be dropped.

Answer: Table S2 contains the metrics of the raw nanopore assemblies while Table 1 contains the metrics of the post-polished assemblies. We kept Table S2 in the supplementary data.

I would prefer that Table S4 be moved out of the supplement and into the main article. Details on polishing effects are important for understanding the ONT platform, and therefore it is unfortunate to bury them in the supplement. I might also argue that Tables S8 and S9 are unfortunate to maroon in the supplement, as these are demonstrating the value of the long read assembly.

Answer: We moved Tables S4, S8 and S9 in the main text.

Reviewer number: 2

Reviewer's report:

Authors present a survey on de novo assembly of yeast genomes using Oxford Nanopore MinION sequencer. Authors assembled a total of 22 yeast strains, the *Saccharomyces cerevisiae* S288C used to assess the quality and performances of the assemblers and data, and 21 strains selected for their diversity and spread.

They compare various types of data that can be produced with MiniION (e.g. 2D and 1D reads) and different MiniION chemistries (R7.3 and the more recent R9).

Data is assembled using 4 different MiniION only assemblers: Canu], Miniasm, SMARTdenovo and ABruijn.

They perform many assemblies with different types of data as input. They also use Illumina read to error correct the final assembly with Pilon.

In general I like the paper, it is a snapshot of the current status of de novo assembly with MiniIon and gives the possibility to a reader to have an idea of what tools to use and what results to expect.

I have some concerns that I want the authors to address:

- they often say in the text "kept the best assembly for each software" (e.g., page 6 line 141). They employ many metrics to discuss about assembly (contiguity, gene coverage, indels) and I like it a lot, but it is not clear how they select the "best" assembly. If for example they choose the best assembly based only on contiguity they might be constantly choosing assemblies affected by many errors, while less contiguous assemblies might be characterized by more correct sequences

Answer: We modified the text at several locations (lines 143,144,150,155,159 and 228) and we replaced "best assembly" by "most contiguous assembly". Indeed, we selected the best assembly for each method based on contiguity metrics (N50, Number of contigs and cumulative size).

- page 6 line 152: "a high proportion of mismatches and indels" : this needs to be more specific, what is "high proportion"?

Answer: We modified the corresponding sentence (lines 153-154), and replaced "a high proportion of mismatches and indels" by "the higher number of mismatches and indels". Furthermore we added a reference to Table S2 which contains the metrics (number of mismatches and indels) of the nanopore-only assemblies before the polishing step.

- The abstract is pretty positive about using only MiniIon data in de novo assembly, or at least that is my impression. Moreover, from the abstract and from the introduction part I was expecting to read about a MiniIon only evaluation and comparison. Instead, the assemblies presented in Table 1 and the various discussions on the evaluation show that all MiniIon assemblies needed Pilon (and therefore the 300X Illumina coverage) to be corrected. Moreover, to finish up the genomes 8Kbp and 20Kbp library have been used, and I assume these are MP Illumina libraries. Therefore, I am now pretty skeptical about the ability of miniion to assemble alone yeast genome... I think that the abstract needs to be toned down and point more the need of complementary technologies to obtain a final assembly.

Answer: We take into account the reviewer comments and change the last sentence of the abstract (lines 56-59) to better reflect current issues with nanopore-only assemblies. All the assemblies were based on nanopore and Illumina paired-end sequencing; we didn't sequence any mate-pairs data. The nanopore-only assemblies show an accurate structure (organization of genomic elements, like genes or transposons) but the final consensus of those assemblies still remains problematic.

I want to point out this paper (I am one of the authors) <https://gigascience.biomedcentral.com/articles/10.1186/s13742-015-0094-1>

In this paper a multi-technology approach is followed combining Illumina, PacBio, and Optical Maps on a yeast genome. In case you have a similar variety to the one assembled in this paper would be nice to compare the assembly presented in the paper with the MiniIon assembly... This is a plus, but I think it would really show potentials of MiniIon.

Answer: It would be very interesting to compare the results of a MinION sequencing strategy and the multi-technology approach (i.e. Illumina, PacBio and Optical Map), which was used to assemble the genome of the Dekkera bruxellensis yeast in Olsen et al. - GigaScience 2015. However, because we sequenced Saccharomyces cerevisiae genomes in our study, this is something impossible. Indeed, Dekkera bruxellensis and Saccharomyces cerevisiae are not closely related species and their genomes are very different in terms of chromosome number and there is no conservation of synteny at all.

Reviewer number: 3

Reviewer's report:

This is an excellent, timely, and well put together study. The results will be greatly helpful to many working on integrating this technology into the genome sequencing ecosystem.

Lines 195-196 describe read polishing with Pilon. It would be helpful to indicate what this depth of coverage was used to polish with the 2x250bp - I realize its in the table legend but could be helpful to include in the text here. also might be helpful to know if 300x is really needed to correct / polish well - would 100x work equally well?

Answer: We modified the text accordingly; we added the Illumina coverage that was used during the polishing step (line 198). Concerning the optimal coverage, we agreed with the reviewer and we performed several polishing with subsets of Illumina PE reads (from 25X to 300X). We added several sentences in the text (lines 205-209) and a supplementary figure (Figure S6) to show that a low coverage (100X or less) is sufficient to correct the consensus of nanopore-only assemblies.

Lines 209-216. Comparing the SPAdes Illumina assembly to the Nanopore only assembly -- The Table presents the QUAST(?) results that gene completeness is actually lower in the Illumina-only assembly but these are mostly indel free? Could be mentioned in the text here?

Answer: We added a sentence to underline these features (lines 214-216).

Doesn't SPAdes also have a option for co-assembly with Illumina + MinION data? Did this produce a useful / comparable assembly ?

Answer: That's right, but we'd like to focus our comparison on nanopore-only assemblers. However, we launched Spades with Illumina and nanopore data. The output assembly was composed of 143 contigs with a N50 near 250Kb. Although the assembly is less fragmented than the Illumina-only assembly it still remains highly fragmented compared to nanopore-only assemblies.

Lines 240 - 257. Sequencing the additional strains. It was unclear how the Pilon polishing is done here - the authors say 300x Illumina paired-end reads - are these from the same strain? Were illumina libraries made and sequenced for each strain or was this using the 1002 genome data ? (the 1002 site says it used 2x102 bp?)

Answer: We modified the text to better describe the pilon process (line 257) and we added a section in the method chapter to explain how the 22 genomes have been assembled (lines 504-509).

One idea I had in reading the manuscript. An additional type of repeat variation that is seen in *Saccharomyces* and other yeasts is the changes regarding simple repeats. These are particularly interesting in context when they fall within context of genic region generating instability that leads to phenotypic variation as the authors I am sure are aware. This was explored through PCR and sequencing in multiple strains by Verstrepen et al Nat Gen 2005 - in particular FLO1 has variable repeated regions easy to pick out. I searched FLO1 against the assemblies and found nice example of expanded repeat in the gene either matching the FLO5 or FLO1 copies. I worked up the example here

https://gist.github.com/hyphaltip/9f5256854f7a049ad81847c4740ece94#file-flo_loci-table

So it looks like there is variability in the size of the repeats in a few of these strains. Up to the authors if this is worth remarking on but it might be something that could also be better resolved than in Illumina assembly.

Answer: This is definitively an interesting comment. Indeed, there are repeated regions in the FLO1 gene, which additionally have an impact on the phenotypic diversity (e.g. adhesion, flocculation or biofilm formation). FLO1 is 4.6 kb long and contains a variable number of repeats of approx. 100 nt, separated by a 45-nt sequence. Consequently, these repeated structures having a small size can be resolved using Illumina sequencing data. In the frame of our study, we really wanted to focus on larger repeated structures, i.e. involving entire genes such as ENA and CUP genes tandem arrays. Indeed, long read sequencing technologies should have a high resolution compared to short read strategies.

Excellent description of methods, versions of software used, and providing reproducible methods. Though it changes rarely, it may be useful to spell out the exact version of the S288C genome assembly and GFF files used in validation.

Answer: We'd like to thank the reviewer for its conscientious reading of the manuscript as well as its suggestions of improvements.