

1 **Long-read only assembly of *Drechmeria coniospora***
2 **genomes reveals widespread chromosome plasticity and**
3 **illustrates the limitations of current nanopore methods.**

4

5 Damien Courtine¹, Jan Provaznik², Jerome Reboul^{1*}, Guillaume Blanc³ Vladimir
6 Benes² and Jonathan J. Ewbank¹⁺

7 ¹Aix Marseille Univ, CNRS, INSERM, CIML, Turing Centre for Living Systems,
8 Marseille, France

9 ²European Molecular Biology Laboratory (EMBL), GeneCore, Heidelberg, Germany.

10 ³Aix Marseille Univ., Université de Toulon, CNRS, IRD, MIO UM 110, 13288, Marseille,
11 France

12 *Current address, Institut de Génétique Moléculaire de Montpellier, Montpellier,
13 France

14 + ewbank@ciml.univ-mrs.fr

15

16 **Abstract**

17 Long read sequencing is increasingly being used to determine eukaryotic genomes.

18 We used nanopore technology to generate chromosome-level assemblies for 3

19 different strains of *Drechmeria coniospora*, a nematophagous fungus used extensively

20 in the study of innate immunity in *Caenorhabditis elegans*. One natural geographical

21 isolate demonstrated high stability over decades, whereas a second isolate, not only

22 had a profoundly altered genome structure, but exhibited extensive instability. We

23 conducted an in-depth analysis of sequence errors within the 3 genomes and

24 established that even with state-of-the-art tools, nanopore methods alone are
25 insufficient to generate eukaryotic genome sequences of sufficient accuracy to merit
26 inclusion in public databases.

27

28 **Background**

29 *Drechmeria coniospora* is an obligate parasitic fungus belonging to the order of
30 Hypocreales. This fungus forms spores that adhere to the cuticle of a range of different
31 nematodes to infect them [1]. We adopted *D. coniospora* strain ATCC-96282, derived
32 from a strain isolated in Sweden, as a model pathogen for *Caenorhabditis elegans* 20
33 years ago [2]. We have cultured this strain, referred to here as Swe1, continuously
34 since then, using it to understand innate immune mechanisms in its nematode host
35 [3,4].

36 As part of our characterization of the interaction between *D. coniospora* and *C.*
37 *elegans*, in 2013, we extracted DNA from our laboratory strain of the time (referred to
38 here as Swe2), and determined its genome. Despite attempts to complete the
39 assembly, the Swe2 genome remained fragmented, with an N50 of 3.86 Mb [5]. In
40 addition to the genome of Swe2, a second *D. coniospora* genome is available (referred
41 to here as Dan2) [6], derived from a strain related to a Danish isolate (Dan1; Figure 1).
42 Although corresponding to a chromosome level assembly, this latter genome still
43 contains large stretches (up to 500 kb) of undetermined sequence. In this study, we
44 used Oxford Nanopore Technology (ONT) long-read sequencing to assemble
45 complete fungal genomes. This revealed that the 2 isolates (Swe1 and Dan1) display
46 strikingly different levels of genomic stability. We provide a detailed analysis that
47 illustrates the continuing challenges to using only ONT long-read sequencing for

48 genome assembly. As the genome sequences were of insufficient quality to allow
49 accurate gene prediction, we polished the genomes using short DNA reads to generate
50 high-quality sequences, providing a resource for future comparative studies.

51

52 **Result**

53 An all-against-all *in silico* genome comparison of the 2 publicly available *D. coniospora*
54 genome sequences, for Dan2 [6] and Swe2 [5], indicated the presence of extensive
55 genomic rearrangements (Figure 2A). These could reflect real differences or assembly
56 errors in one or both genomes. We directly confirmed one major rearrangement by
57 PCR (Figure 2B, C), suggesting that the differences could be real. To characterise this
58 genomic plasticity, we determined the genomes of 3 strains related to the 2 that had
59 been sequenced previously (Figure 1). We used ONT nanopore sequencing to
60 generate long reads and current assembly tools to construct chromosome level
61 assemblies for all 3 strains (Supplementary Fig. S1, Supplementary Table S1). Manual
62 curation allowed complete ca. 30 kb mitochondrial genomes to be predicted from the
63 assemblies generated by Canu [7].

64 All 3 nuclear genomes were divided in 3 similarly sized chromosomes, an unusual
65 arrangement for such a fungus, as previously noted by Zhang *et al.* for Dan2 [6]. For
66 the 2 strains related to Swe2, there was almost complete synteny of their nuclear
67 genomes. Inspection of the one anomalous region in Swe1 where synteny broke down
68 revealed that it was supported by only one long (215 kb) read and corresponded to a
69 local discontinuity in the read coverage, as well as a break in the alignment between
70 Canu-generated contigs and unitigs. All these factors indicated that this was an
71 assembly artefact with a contig misassembled on the basis of an individual very long

72 chimeric read (Supplementary Fig. S2). The same was true for the distinct unique non-
73 syntenic region of the Swe3 assembly (Supplementary Fig. S3).
74 These were exceptional cases since the overwhelming majority of chimeric reads were
75 identified and either trimmed or excluded from the assembly process by Canu
76 (Supplementary Fig. S4, Supplementary Fig. S5). An in-depth analysis of the Swe1
77 chimeric reads revealed that a large proportion was in fact the consequence of
78 sequencing errors. In almost 40% of cases (1010 / 2566), the two regions flanking the
79 presumptive site of chimerism mapped to within 50 nucleotides of each other on the
80 corresponding single scaffold. There was no discernible pattern to the distribution of
81 this interval in the remaining candidate chimeric reads (Supplementary Fig. S6A-B),
82 nor where there any regions that were more likely to be the site of chimeric junctions
83 (Supplementary Fig. S6C).
84 Notably the single chimeric read that escaped censoring, leading to a misassembly of
85 Swe1, was not identified by the dedicated tool YACRD, but was flagged as anomalous
86 in reads recalled by Guppy (see Methods). This is an indication of the continuing
87 improvement to base-calling tools. Also, these specific Swe1 and Swe3 misassemblies
88 were absent from the corresponding chromosome assemblies produced by the *de*
89 *novo* assembler Flye [8] (Figure 3A). This latter, however, introduced other assembly
90 artefacts, including an erroneous fusion of contigs for the Dan1 assembly. This could
91 not be ascribed to the inclusion of chimeric reads, but rather appeared to result from
92 the incorrect treatment of repeat sequences, including telomeric repeats at the
93 extremity of one of the fused contigs (Figure 3B-D). These results illustrate the interest
94 of using more than one tool to aid in genome assembly. Therefore, starting with the

95 Canu-generated sequences, we manually corrected anomalous regions and thereby
96 produced assemblies for Swe1 and Swe3 that were entirely collinear (Figure 4A).
97 These 2 genomes have 3 large chromosomes (8.5 Mb, 11.6 Mb, 11.6 Mb), each with
98 identifiable telomeric [9] and centromeric regions, indicating that the overall genome
99 structure has remained constant over 20 years of laboratory culture. This allowed us
100 then to use the Swe1 sequence to scaffold the fragmented Swe2 genome (Figure 4B).
101 To our great satisfaction, we were able to produce an entirely collinear chromosome
102 scale assembly. Thus, it appears that there were no assembly errors in the published
103 Swe2 genome, it was simply incompletely scaffolded. This applies equally to the
104 genomic regions containing copies of some mitochondrial genes that we previous
105 suggested might indicate assembly errors [5]. They were revealed to be accurate; *D.*
106 *coniospora* has nuclear paralogous copies of 10 mitochondrial protein-coding and 15
107 tRNA genes (so called *numts* sequences [10]). These results give further support to
108 the existence of long-term stability of the genome of the Swe2 related strains. A whole
109 genome comparison between Swe1 and Dan2, however, revealed multiple and
110 extensive genome rearrangements, involving intra- and inter-chromosomal
111 translocations and inversions (Figure 4C).
112 Using the same strategy described above, we assembled and polished the Dan1
113 genome to give chromosome-level sequences. When we compared Dan1 and Dan2,
114 we were surprised to find 2 major events of reciprocal exchange of chromosome ends,
115 and an intra-chromosomal inversion (Figure 4D). These events were supported in a
116 coherent and consistent manner by all the available data (Supplementary Fig. S7). In
117 other fungal species, such chromosomal rearrangements have been reported to be the
118 result of ectopic recombination between non-allelic homologous sequences, including

119 repeated DNA elements [11,12]. A search of the 50 kb regions flanking each break
120 point for transposable elements [13] and repetitive DNA families [14], failed to reveal
121 any significant repeat sequence signature (see Supplementary Methods). As the Dan2
122 assembly is of high confidence, supported by long reads and optical mapping [6], given
123 the short time of *in vitro* culture that separates it from Dan1, this suggests that the
124 genome of the Dan1 isolate is not stable.

125 In alignments of the sequence of Swe1, generated using only nanopore reads, with
126 that of Swe2, there were stretches of complete nucleotide identity extending over more
127 than 25 kb. This is a testament to the general reliability of nanopore sequencing. We
128 therefore identified the complete set of proteins identical in Swe2 and Dan2
129 corresponding to single copy, single exon genes (see Methods). These would be
130 expected to be present in the newly assembled Swe1, Swe3 and Dan1 genomes.
131 Indeed, using these 305 genes as a query, we could identify homologous sequences
132 for each in all 3 genomes. Less than 1/6 of the corresponding genes, however, were
133 predicted to encode full-length proteins in any of the 3 new genomes (Figure 5A). While
134 nanopore reads are very useful for genome assembly, they suffer from a high error
135 rate, especially in homopolymer stretches. Sequence quality can be improved using
136 polishing tools that aim to ameliorate consensus sequences generally by going back
137 to raw reads and applying integrative algorithms [15]. In our case, applying current best
138 practices, while providing a very substantial improvement (up to 5-fold in the best
139 case), did not take the prediction level beyond 82% accuracy. The quality of the
140 prediction seen with the Dan1 genome was strikingly lower than the other 2 genomes
141 (Figure 5A, Supplementary Table S1).

142 Inspection suggested that the majority of errors were in homopolymer sequences, as
143 expected, with nucleotide insertions and deletions leading to alterations of the reading
144 frame. To investigate this poor homopolymer predictive performance systematically,
145 we computed the number of G/C or A/T homopolymer stretches of at least 4
146 nucleotides for each of the 305 genes. We plotted these values, indicating the
147 proportion of genes that encoded the expected full-length predicted protein for each of
148 the 3 genomes. While there was the expected inverse relationship between accuracy
149 and the number of homopolymer stretches, there were striking exceptions. Curiously
150 some of these exceptions were specific to a single genome (Figure 5B-D). Further,
151 and unexpectedly, polishing introduced more nucleotide insertion errors than
152 deletions, frequently on the basis of tenuous read support. Overall, however, there was
153 no obvious pattern to explain why errors were introduced, given the underlying reads
154 used to build the consensus sequence (Supplementary Table S2).

155 During the inspection of the assembled and polished genomes, we found two other
156 types of anomalies. The first concerned the regions flanking the nuclear genomic
157 copies of mitochondrial genes (*numts*), where polishing added short extraneous low
158 complexity sequences (average length 15 nt, mainly As or Ts), for which, surprisingly
159 there was no sequence support from the reads used by the assembler (Figure 6A).
160 This probably arose because of the very high nucleotide similarity between regions of
161 the nuclear and mitochondrial genomes that extended across more than 25 kb,
162 including a repeat of 9.8 kb (Supplementary Fig. S8A-B). Notably, despite using high
163 coverage ONT long reads, we could not establish with absolute certainty the precise
164 copy number for the unit sequence in the Swe genomes (Supplementary Fig. S8B).

165 In the second case, for the Swe3 genome, a large (ca. 10 kb) region, with a complex
166 sequence, well supported by the Canu corrected and trimmed reads, was inexplicably
167 excluded from the initial Canu assembly and only imprecisely restored by polishing
168 (Figure 6B-C, Supplementary Fig S8C). Here, while there was no evidence for
169 repeated DNA elements on both sides of the point of sequence discontinuity, there
170 was a single such 1.2 kb duplication (Supplementary Fig S8C-D). These few regions
171 were identified because of discontinuities in the depth of read coverage, which
172 otherwise was remarkably constant across the complete genomes. With the resolution
173 of these assembly errors, we were able to generate complete genomes of high overall
174 structural quality using ONT long reads only.

175 As explained above, however, these assemblies were not of sufficient sequence
176 quality to allow accurate gene prediction. Therefore, to extend our analysis, we used
177 Illumina sequencing to generate very deep short read coverage for the Swe1, Swe3
178 and Dan1 genomes. This allowed high quality final sequences to be generated for all
179 3 strains. While short-read-based polishing did not alter the global structure, it allowed
180 homopolymer length errors to be corrected and the generation of entirely contiguous
181 chromosome sequences (Supplementary Table S1).

182 To confirm the correctness of the short-read polished assemblies, we returned to our
183 305 single copy orthologues. After the short-read polishing, all 305 genes could be
184 identified in each of the 3 genomes (Supplementary Table S1). We also benchmarked
185 our successive assemblies using BUSCO that searches for a set of universal single-
186 copy orthologues (USCO) by sequence similarity. While the initial genome assemblies
187 gave low scores, with roughly 65% of complete USCOs and 35% fragmented or
188 missing (Table 1), after long-read polishing the score for complete USCOs increased

189 up to as high as 97%. Given the demonstrably low quality of the genome sequences
190 (Figure 5), we investigated the basis of this disparity. We identified among the USCOs
191 those that corresponded to single exon genes in the Dan2 and Swe2 reference
192 genomes. These genes were then used as queries for high-stringency searches of the
193 Dan1, Swe1 and Swe3 genomes at successive steps of assembly and polishing and
194 the results compared to the results of the corresponding BUSCO analysis. While
195 BUSCO gave no false negatives, it gave a large number of false positives, except in
196 the analysis of the short-read polished genomes (Figure 7A). These arose because
197 BUSCO was not sufficiently sensitive to the presence of short indels. As an example,
198 the Swe1 gene corresponding to RJ55_06485 had the expected sequence after short-
199 read polishing. Two errors in homopolymer sequences led to 2 frameshifts in the
200 unpolished assembly. One of these was corrected by long-read polishing, but for the
201 other there was an over-compensation, leading to a different frameshift (Figure 7B). In
202 both assemblies, these errors were compatible with open-reading frames that
203 collectively reconstituted a close ortholog of RJ55_06485 leading to the erroneous
204 BUSCO result. As discussed below, this analysis highlights the fact that BUSCO
205 scores based on sequence alignments are not an appropriate measure for ONT-only
206 eukaryotic genomes. The BUSCO score rose to nearly 99% after the short-read
207 polishing. In this case, the figures accurately reflect genome completeness and quality
208 (Figure 7A). These figures are comparable to those for the previous Dan2 and Swe2
209 assemblies. The new Swe1, Swe3 and Dan1 genomes therefore represent the starting
210 point for future detailed analysis to characterise the molecular evolution of *D.*
211 *coniospora*.

212

213 Table 1: BUSCO results

Strain	Assembly	Complete	Complete: single	Complete: duplicated	Fragmented	Missing
Dan1	Canu curated	820 (62.4%)	820 (62.4%)	0 (0%)	259 (19.7%)	236 (17.9%)
Dan1	Long-read polish	1187 (90.3%)	1187 (90.3%)	0 (0%)	62 (4.7%)	66 (5%)
Dan1	Short-read polish	1297 (98.6%)	1296 (98.6%)	1 (0.1%)	9 (0.7%)	9 (0.7%)
Dan2	[6]	1298 (98.7%)	1297 (98.6%)	1 (0.1%)	8 (0.6%)	9 (0.7%)
Swe1	Canu curated	869 (66.1%)	868 (66%)	1 (0.1%)	243 (18.5%)	203 (15.4%)
Swe1	Long-read polish	1266 (96.6%)	1266 (96.6%)	0 (0%)	21 (1.6%)	28 (2.1%)
Swe1	Short-read polish	1296 (98.6%)	1295 (98.5%)	1 (0.1%)	8 (0.6%)	11 (0.8%)
Swe2	[5]	1296 (98.6%)	1294 (98.4%)	2 (0.2%)	9 (0.7%)	10 (0.8%)
Swe3	Canu curated	859 (65.3%)	858 (65.2%)	1 (0.1%)	243 (18.5%)	213 (16.2%)
Swe3	Long-read polish	1274 (96.9%)	1274 (96.9%)	0 (0%)	17 (1.3%)	24 (1.8%)
Swe3	Short-read polish	1295 (98.5%)	1294 (98.4%)	1 (0.1%)	9 (0.7%)	11 (0.8%)

214 Percentage of each category of the expected 1315 USCOs for different genome
 215 assemblies. Of the 11 USCOs missing in Swe1 and Swe3, 10 are also absent from
 216 Swe2, and 9 from Dan1 (and Dan2). These are therefore likely to be real gene losses
 217 in *D. coniospora*, so that only 2 USCOs (0.2%) at most are missing.

218

219 Discussion and conclusion

220 Previous genome assemblies for *D. coniospora* required a combination of sequencing
 221 approaches [5,6]. Here, using only long reads and Canu, we produced the first
 222 complete circular mitochondrial genome for *D. coniospora* and were able to generate
 223 chromosome-scale assemblies for the nuclear genome. The rare misassembled
 224 contigs, formed by Canu because of single very long chimeric reads, as previously
 225 described [16], could be detected by read coverage anomalies and comparisons with
 226 unitigs, suggesting that solutions to avoid their creation could be implemented within
 227 Canu. The majority of reads that were flagged as chimeric arose from sequencing or
 228 polishing errors. They reflected a short (<50 bp) discrepancy between the individual
 229 reads and the final sequence. There was no indication of any sequence bias at the

230 break points of the remaining chimeric reads supporting the notion that these reads
231 arise from too rapid reloading of the sequencing pore [17].

232 The use of other genome assembly tools, and the comparison of assembly
233 discrepancies is an additional method to produce high confidence genomes. Here, we
234 used Flye that for these genomes required run times that were ten-fold shorter than
235 Canu. A comparison of the assemblies highlighted ambiguous regions in the genome
236 that could then be resolved by manual inspection. On the other hand, Flye was
237 confounded by telomeric repeats. Since telomeres can be identified on the basis of
238 their sequence, there is also clear room for algorithmic improvement to Flye through
239 the explicit definition of chromosome ends.

240 One clear and well-established advantage of using long reads is the possibility of
241 resolving very extended stretches of complex tandem repeats (VeCTRs) [18] and other
242 repetitive sequences including centromeres. These correspond to most of the breaks
243 in the continuity of the published Swe2 genome. In addition to acrocentric regional
244 centromeres, Zhang *et al.* reported the presence of a vestigial centromere from a
245 putative chromosomal fusion event [6]. These were also found in the fully assembled
246 Swe1 and Swe3 genomes, indicating that chromosomal fusions were present in the
247 common ancestor of the Swe1 and Dan1 strains.

248 For Swe1, Swe3 and Dan1 we were able to reconstruct complete mitochondrial
249 genomes, with features typical of fungi of the order Hypocreales. On the other hand,
250 unlike Dan1 (and Dan2), the nuclear genomes of Swe1 and its derivatives Swe2 and
251 Swe3, contained different numbers of copies of sequence very similar to parts of their
252 own mtDNA. This type of event, and more generally repeated regions with long and

253 nearly identical sequences are more readily detectable with long reads [19], and are
254 particularly challenging for polishing even with short reads [20].

255 The duplication of mitochondrial genes in the nuclear genome has been described in
256 other fungal genomes [10] and must have occurred after the divergence of Dan1 and
257 Swe1. Despite this genome plasticity, even after 20 years of continuous laboratory
258 culture the Swe1 and Swe3 genomes were entirely collinear. This contrasts with the
259 rearrangements seen between the Dan1 and Dan2 genomes that in principal should
260 be from strains that have had little opportunity to diverge (L. Castrillo, Curator, ARS
261 Collection of Entomopathogenic Fungal Cultures, personal communication). It will be
262 interesting in the future to characterize the reasons for the marked difference in
263 genomic stability between Dan1 and Swe1.

264 The accuracy of ONT long read sequencing is increasing because of improvements in
265 the chemistry used, signal detection, as well as base-calling [21]. Despite good read
266 depth, however, our assemblies were not of sufficient quality at the nucleotide level to
267 allow accurate gene prediction. Further, we noted that although polishing using only
268 long reads dramatically increased overall sequence accuracy, it introduced errors
269 around the *numts*. Similar errors during polishing of near identical sequences has been
270 noted in ONT-based metagenomic studies [22]. Despite these limitations, research
271 groups are publishing and submitting to public sequence databases genomes for fungi,
272 plants and animals based on nanopore sequencing alone (86 for Eukaryotes in
273 addition to the 134 Bacterial genomes in “Assembly” from GenBank release 236 from
274 the 2020/02/15). This is problematic as low-quality genome sequences compromise
275 the accuracy of sequence similarity searches in public databases. On the basis of our
276 results, a re-analysis of the completeness of these “nanopore-only” genomes is

277 merited, to confirm that they are indeed low quality. Similar concerns do not apply to
278 fungal genomes assembled using only long reads generated with Pacific Bioscience
279 technology [23] as these do not suffer from the intrinsic problem of homopolymer length
280 errors that we found to be the most significant quality barrier when using ONT reads.
281 On the basis of our detailed analysis and in line with the consensus regarding *de novo*
282 assembly with ONT long reads (e.g. [24]), we polished our 3 assemblies with short
283 reads. This greatly improved their quality.

284 Regarding the homopolymer sequence errors, as noted above, they were not
285 consistent across the sequenced genomes; even between Swe1 and Swe3 there were
286 instances of widely differing rates of errors in orthologous genes, despite very similar
287 underlying reads. Indeed there was no clear pattern in the inaccuracies, which will
288 render bioinformatics approaches to remedy this problem more difficult. On the other
289 hand, the errors were more often over-prediction of homopolymer length, despite
290 having a majority of reads supporting the correct sequence. It is possible that polishing
291 tools have not kept pace with improvements in base-calling, leading to an over-
292 compensation in the inference of homopolymer length.

293 It is standard practice to check the completeness of *de novo* genome assemblies with
294 a strategy based on the detection of predicted groups of conserved orthologous
295 proteins. One popular and much cited tool is BUSCO [25] which was developed before
296 ONT-based sequencing became prevalent. Since BUSCO relies on *in silico* translation,
297 small indels can be overlooked as the resulting virtual sequence can be recapitulated
298 despite a frameshift. This explains the disparity between the BUSCO results and our
299 own analyses that were deliberately restricted to mono-exonic genes. Contrary to
300 BUSCO, our analysis indicated that about 1/5 of the genes after long-read polishing

301 had an incorrect sequence. Current BUSCO-type approaches, based on sequence
302 similarity and not excluding genes with improbably short introns, cannot be used as a
303 quality metric for ONT-only assemblies, and are appropriate only after short-read
304 correction.

305 In conclusion, nanopore long read sequencing provides a powerful way to assemble
306 complex genomes with limited manual curation but still fall short of the quality required
307 to produce publishable eukaryotic genomes. In our case, it has revealed new
308 information about genome plasticity in *D. coniospora* and provided a backbone that will
309 permit future detailed study to characterize gene evolution in this important model
310 fungal pathogen.

311

312

313 **Methods**

314 DNA extraction:

315 *D. coniospora* spores were cultured in liquid NGMY medium [26] at 37°C for 5 days.
316 Fungal DNA was extracted according to a published protocol (from p13 onwards of
317 [27]) [28], with the following modifications: instead of centrifugation to collect DNA after
318 precipitation with isopropanol, we recovered the DNA filaments with a glass hook,
319 washed and dried them as described [29] and resuspended the DNA without agitation
320 in Tris-EDTA buffer.

321

322 Nanopore sequencing library preparation:

323 Libraries were prepared for sequencing on GridION with the ligation sequencing kit
324 SQK-LSK109. The GridION sequencing was run on flowcell FLO-MIN106 for 47, 48
325 and 48 hours using 972, 660 and 610 ng of DNA (for Swe3, Swe1, Dan1 respectively)
326 and MinKNOW 2.1 v18.06.2.

327

328 Illumina sequencing library preparations:

329 The same DNA samples were used to prepare paired-end libraries with insert size of
330 circa 680 bp, following the manufacturer's instructions for the kit NEBNext® Ultra™ II
331 DNA (New England Biolabs Inc. Ipswich, MA, USA). The libraries were sequenced
332 using an Illumina NextSeq500 system (s/n NB501764).

333

334 Basecalling, adaptor trimming and chimeric read detection:

335 For a first assembly, reads were basecalled at the EMBL using Guppy v1.5.1 (Oxford
336 Nanopore Technologies). For subsequent polishing, we used Guppy v3.0.3 (with

337 parameters -c dna_r9.4.1_450bps_hac.cfg), then adaptors were trimmed with
338 Porechop v0.2.4 [30] with default parameters. YACRD v0.5.1 [31] with the
339 subcommand chimeric and the option --filter was used to remove chimeric reads.

340

341 Whole genome alignments:

342 Genomes were aligned using LAST v979 [32]. A database was first generated (last-db
343 -cR01), and then lastal and last-dotplot with default parameters were used to generate
344 respectively an alignment file and a dot-plot. For the circular visualization of genome
345 alignments, we used the command lastal with -f BlastTab parameter, then parsed the
346 alignment to filter out short alignments and generate the links file needed by Circos
347 [33].

348

349 Mapping of long reads:

350 Validation of genomes during and after assembly involved rounds of read mapping.
351 Reads were aligned with Minimap2 v2.16r922 [34] (with parameters -ax map-ont). The
352 resulting mapping file was processed with Samtools v1.9 [35] to obtain a sorted BAM
353 file (samtools view -bS -q 1 -F 4; samtools sort; samtools index). Mapping results were
354 visualized with IGV v2.5.0 [36].

355

356 Genome assembly:

357 Assemblies were performed with Canu v1.7 [37] and the parameters useGrid=False,
358 genomeSize=30m, correctedErrorRate=0.16 with reads basecalled by Guppy v1.5.1.
359 For the manual curation of the assemblies, we generated whole assembly alignments
360 and dot-plots of Swe1, Swe2 and Swe3 two by two. For Swe1 and Swe3, Canu contigs

361 were ordered by synthesizing the results from the 3 possible all-against-all alignments.
362 To confirm a link between two contigs, we employed the following strategy: when a
363 contig of the Swe1 assembly spanned two contigs of Swe3, long reads of Swe1 present
364 in this spanning area were extracted from the Swe1 corrected and trimmed reads
365 provided by Canu. Then this set of reads was mapped on Swe2 and Swe3 assemblies.
366 The two targeted contigs of Swe3 were considered 'linked' if different parts of several
367 unique reads mapped on the two Swe3 contigs ends. If the reads that supported the
368 link had different mapping orientation (forward or reverse), one contig was
369 complemented before the last step (see *Solving links between contigs*) to ensure a
370 correct orientation of the final chromosome.

371 To guide correct assembly, we also searched for centromeres in the contigs. They
372 were identified as highly duplicated regions in the all-against-all alignment dot-plots
373 produced by LAST. The identification of the repeated canonical telomeric sequence
374 (TTAGGG)_n [9] and its reverse complement (CCCTAA)_n at the beginning or end of
375 certain contigs allowed the identification of chromosome ends. The Dan1 assembly
376 was manually curated using a similar strategy with the Dan2 genome as a reference.

377

378 Solving links between contigs:

379 Overlaps between linked contigs were identified by a BLASTN [38] alignment of their
380 last 100 kb. Any duplicate sequence was trimmed out from one contig and both contigs
381 were joined. The inferred junction was then validated by verification of the underlying
382 read support. For the linked contigs that did not overlap, the sequence in the gap was
383 extrapolated from the reads that matched and extended the ends of contigs, on the
384 basis of alignments at the last 1 kb of each contig. These sequences were aligned with

385 MAFFT v7.427 [39]. The alignment was visualized with SeaView [40], and only the
386 portion of the alignment strictly between the two contigs sequences was kept. Seaview
387 also generated a consensus sequence (on the basis of 60 % sequence identify by
388 default). The resulting sequence was inserted between the two contigs to link them
389 and the supposed continuity verified by a further cycle of read mapping.

390

391 Assembly polishing with long-reads:

392 Genome polishing was carried out with 2 or 4 iterative runs of Racon v1.4.2 [41] and
393 parameters -m 8 -x -6 -g -8 -w 500, and a run of Medaka v0.8.1 (Oxford Nanopore
394 Technologies) with the parameter -m r941_min_high.

395

396 Mitochondrial genome circularization:

397 Canu assembles small circular elements as contigs with tandem duplications of the
398 element. We resolved the mitochondrial genomes as recommended by Canu's authors
399 [7]. MUMmer suite v4.0.0.beta2 [42] was used to align the contig identified as the
400 putative mitochondria on itself with NUCmer and parameters --maxmatch --nosimplify.
401 Coordinates of a full copy were identified with the show-coords command and -lrcT
402 parameters.

403

404 PCR:

405 PCR was carried out to test a genome rearrangement between Swe2 and Dan2
406 genomes, with primers P1F (GAGATATCGAACGTCGCATGG), P1R
407 (ACATCAAGCCTTTGTCGAGGA), and P3F (GCTCAGGACCGACGTACAAG). PCR
408 reactions were run according to the GoTaq® G2 Flexi DNA polymerase instructions

409 (Promega), with 50 ng of template DNA, 1 mM of each forward and reverse primer, in
410 a final volume of 25 μ L. The reaction started by initial denaturation at 95°C for 2 min,
411 followed by 30 amplification cycles (95°C for 30 sec, 60°C for 30 sec and 72°C for 30
412 sec), and a final elongation for 5 min at 72°C.

413

414 Defining a set of 305 identical proteins:

415 Identical proteins shared by the two *D. coniospora* genomes available (Swe2 and
416 Dan2) were recovered using a reciprocal best BLAST [38] hit strategy on the two
417 proteomes. Proteins that were duplicated in one or both genomes were filtered out.
418 The set was further refined by only retaining proteins corresponding to mono-exonic
419 genes.

420

421 Assessment of gene sequence in ONT-only assemblies:

422 TBLASTN searches were run using the amino-acid sequence of the set of 305 identical
423 proteins against the different nanopore only assemblies. A gene was considered as
424 correct if the query coverage, *i.e* the ratio of alignment length over the query length,
425 was equal to 1.

426

427 Short read polishing:

428 Short reads were trimmed using Trimmomatic v0.39 [43] with the parameters
429 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:30 MINLEN:36. Then, only paired
430 reads were mapped on assemblies with bwa v0.7.17 [44] and default parameters (bwa
431 index, then bwa mem). The resulting mapping file was converted in BAM, sorted and
432 indexed with samtools. This latter file was used to polish the assembly with Pilon v1.23

433 [45] with the parameters --fix bases --vcf --mindepth 10 --minmq 20 --minqual 15 --
434 changes --diploid. Several iterations were conducted for each strain, until the number
435 of changes was less than 5.

436

437 Flye assembly:

438 An additional *de novo* assembly was performed with Flye v2.4.2 [8], and the parameter
439 --genome-size 32m, using the ONT reads recalled by Guppy v3.0.3.

440

441 Assessing the genome integrity:

442 The genome integrity was assessed with BUSCO v3.1.0 and the curated set
443 *ascomycota_odb9* version 2016-02-13 [25]. A BLASTP search enabled Swe2
444 monoexonic genes present among USCOs to be identified. This list of 219 Swe2 genes
445 was then used as a TBLASTN query against the different assemblies of Swe1 and
446 Swe3. A gene was considered correct when it matched the corresponding Swe2 gene
447 perfectly in length. An analogous analysis was carried out for Dan1, on the basis of the
448 273 Dan2 monoexonic genes that are USCOs.

449

450 Characterisation of chimeric reads:

451 Swe1 reads identified as chimeric by YACRD were aligned on the final (short-read
452 polished) Swe1 assembly. The main alignment was identified using samtools view -F
453 2308. The CIGAR string was then parsed to determine whether the longest residual
454 part of the read was 5' or 3' to the main alignment, thereby giving an orientation to the
455 putative chimeric read and localising the potential chimeric break point. The 500 bp of
456 sequence 5' and 3' of this point were extracted and individually mapped back on the

457 Swe1 final assembly and the number of unique reads in a 10 kb non-overlapping sliding
458 window was calculated. For the reads for which both 500 bp fragments mapped on the
459 same chromosome, the smallest distance between the two fragments was calculated.

460

461 **Data availability**

462 Genomes of the strains Swe1, Swe3 and Dan1 are available on our institute website
463 (<http://www.ciml.univ-mrs.fr/applications/DC/Genome.htm>). All the reads used in this
464 work can be found at the European Nucleotide Archive (ENA) under the study numbers
465 PRJEB35969, PRJEB35970 and PRJEB35971. The raw signal runs are available
466 under the accessions ERR3774158, ERR3774162 and ERR3774163; the FASTQ files
467 of basecalled reads (Guppy v3.0.3) are available under the accessions ERR3997391,
468 ERR3997394 and ERR3997483; the FASTQ files of Illumina paired-end reads are
469 available under the accessions ERR3997389, ERR3997392, ERR3997395. Accession
470 numbers are given in the order Swe1, Swe3 and Dan1.

471

472 **Funding**

473 Supported by institutional grants from the Institut national de la santé et de la recherche
474 médicale, Centre National de la Recherche Scientifique and Aix-Marseille University
475 to the CIML, and the Agence Nationale de la Recherche program grant (ANR-16-
476 CE15-0001-01), and "Investissements d'Avenir" ANR-11-LABX-0054 (Labex
477 INFORM), ANR-16-CONV-0001 and ANR-11-IDEX-0001-02, and funding from the
478 Excellence Initiative of Aix-Marseille University - A*MIDEX.

479

480 **Acknowledgments**

481 The authors thank Yuquan Xu and Liwen Zhang for providing access to the raw
482 sequencing data for Dan2, Xing Zhang for help in preparing DNA samples, Lionel
483 Spinelli for informatic support and Nathalie Pujol and Laurent Tichit for comments.

484

485 **Supplementary data**

486 *supplementary_figures.pdf* contains 8 supplementary figures.

487 *supplementary_table_1.xls* contains the read coverage of the genomes. This table also
488 contains the results for the TBLASTN on the 305 candidate identical proteins.

489 *supplementary_table_2.xls* is a table recording read support, in Swe1, Swe3 and Dan1
490 assemblies, for the predicted correct sequence for each homopolymer stretch in the
491 genes corresponding to 10 protein of the 305 candidate identical proteins.

492 *supplementary_methods.pdf* contains additional methodological details.

493

494 References

- 495 1. Jansson HB, Jeyaprakash A, Zuckerman BM. Differential Adhesion and Infection of
496 Nematodes by the Endoparasitic Fungus *Meria coniospora* (Deuteromycetes). *Appl Environ*
497 *Microbiol.* 1985;49:552–5.
- 498 2. Pujol N, Link EM, Liu LX, Kurz CL, Alloing G, Tan MW, et al. A reverse genetic analysis
499 of components of the Toll signalling pathway in *Caenorhabditis elegans*. *Curr Biol.*
500 2001;11:809–21.
- 501 3. Dierking K, Polanowska J, Omi S, Engelmann I, Gut M, Lembo F, et al. Unusual regulation
502 of a STAT protein by an SLC6 family transporter in *C. elegans* epidermal innate immunity.
503 *Cell Host Microbe.* 2011;9:425–35.
- 504 4. Labeed SA, Omi S, Gut M, Ewbank JJ, Pujol N. The pseudokinase NIPI-4 is a novel regulator
505 of antimicrobial peptide gene expression. *PLoS One.* 2012;7:e33887.
- 506 5. Lebrigand K, He LD, Thakur N, Arguel M-J, Polanowska J, Henrissat B, et al. Comparative
507 Genomic Analysis of *Drechmeria coniospora* Reveals Core and Specific Genetic Requirements
508 for Fungal Endoparasitism of Nematodes. *PLoS Genet.* 2016;12:e1006017.
- 509 6. Zhang L, Zhou Z, Guo Q, Fokkens L, Miskei M, Pócsi I, et al. Insights into Adaptations to a
510 Near-Obligate Nematode Endoparasitic Lifestyle from the Finished Genome of *Drechmeria*
511 *coniospora*. *Sci Rep.* 2016;6:23122.
- 512 7. Canu 1.8 documentation. [https://canu.readthedocs.io/en/latest/faq.html#my-circular-](https://canu.readthedocs.io/en/latest/faq.html#my-circular-element-is-duplicated-has-overlap)
513 [element-is-duplicated-has-overlap](https://canu.readthedocs.io/en/latest/faq.html#my-circular-element-is-duplicated-has-overlap). Accessed 15 Nov 2019
- 514 8. Kolmogorov M. Fast and accurate de novo assembler for single molecule sequencing reads:
515 fenderglass/Flye. <https://github.com/fenderglass/Flye>. Accessed 3 June 2019
- 516 9. Schechtman MG. Characterization of telomere DNA from *Neurospora crassa*. *Gene.*
517 1990;88:159–65.
- 518 10. Hazkani-Covo E, Zeller RM, Martin W. Molecular Poltergeists: Mitochondrial DNA
519 Copies (numts) in Sequenced Nuclear Genomes. *PLoS Genet.* 2010;6:e1000834.
- 520 11. Argueso JL, Westmoreland J, Mieczkowski PA, Gawel M, Petes TD, Resnick MA. Double-
521 strand breaks associated with repetitive DNA can reshape the genome. *Proc Natl Acad Sci U S*
522 *A.* 2008;105:11845–50.
- 523 12. Sun S, Yadav V, Billmyre RB, Cuomo CA, Nowrousian M, Wang L, et al. Fungal genome
524 and mating system transitions facilitated by chromosomal translocations involving
525 intercentromeric recombination. *PLoS Biol.* 2017;15:e2002527.
- 526 13. TransposonPSI: An Application of PSI-Blast to Mine (Retro-)Transposon ORF
527 Homologies. <http://transposonpsi.sourceforge.net>. Accessed 1 Apr 2020
- 528 14. Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, et al. The Dfam database of
529 repetitive DNA families. *Nucleic Acids Res. Oxford Academic;* 2016;44:D81–9.

- 530 15. Senol Cali D, Kim JS, Ghose S, Alkan C, Mutlu O. Nanopore sequencing technology and
531 tools for genome assembly: computational analysis of the current state, bottlenecks and future
532 directions. *Brief Bioinform.* 2018;20:1542–59.
- 533 16. Scheunert A, Dorfner M, Lingl T, Oberprieler C. Can we use it? On the utility of de novo
534 and reference-based assembly of Nanopore data for plant plastome sequencing. *PLoS One.*
535 2020;15:e0226234.
- 536 17. White R, Pellefigues C, Ronchese F, Lamiable O, Eccles D. Investigation of chimeric reads
537 using the MinION. *F1000Res.* 2017;6:631.
- 538 18. Eccles D, Chandler J, Camberis M, Henrissat B, Koren S, Le Gros G, et al. De novo
539 assembly of the complex genome of *Nippostrongylus brasiliensis* using MinION long reads.
540 *BMC Biol.* 2018;16:6.
- 541 19. Schmid M, Frei D, Patrignani A, Schlapbach R, Frey JE, Remus-Emsermann MNP, et al.
542 Pushing the limits of de novo genome assembly for complex prokaryotic genomes harboring
543 very long, near identical repeats. *Nucleic Acids Res.* 2018;46:8953–65.
- 544 20. Watson M, Warr A. Errors in long-read assemblies can critically affect protein prediction.
545 *Nature Biotechnology.* 2019;37:124–6.
- 546 21. Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for Oxford
547 Nanopore sequencing. *Genome Biol.* 2019;20:129.
- 548 22. Somerville V, Lutz S, Schmid M, Frei D, Moser A, Irmeler S, et al. Long-read based de novo
549 assembly of low-complexity metagenome samples results in finished genomes and reveals
550 insights into strain diversity and an active phage system. *BMC Microbiol.* 2019;19:143.
- 551 23. Dal Molin A, Minio A, Griggio F, Delledonne M, Infantino A, Aragona M. The genome
552 assembly of the fungal pathogen *Pyrenochaeta lycopersici* from Single-Molecule Real-Time
553 sequencing sheds new light on its biological complexity. *PLoS One.* 2018;13:e0200217
- 554 24. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and
555 assembly of a human genome with ultra-long reads. *Nat Biotechnol.* 2018;36:338–45.
- 556 25. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
557 assessing genome assembly and annotation completeness with single-copy orthologs.
558 *Bioinformatics.* 2015;31:3210–2.
- 559 26. He LD, Ewbank JJ. Polyethylene Glycol-mediated Transformation of *Drechmeria*
560 *coniospora*. *Bio-Protoc.* 2017;7:e2157.
- 561 27. Fungal DNA extraction protocol.
562 [https://www.pnas.org/content/pnas/suppl/2018/01/08/1715954115.DCSupplemental/pnas.171](https://www.pnas.org/content/pnas/suppl/2018/01/08/1715954115.DCSupplemental/pnas.1715954115.sapp.pdf)
563 [5954115.sapp.pdf](https://www.pnas.org/content/pnas/suppl/2018/01/08/1715954115.DCSupplemental/pnas.1715954115.sapp.pdf). Accessed 16 Apr 2020.
- 564 28. Kjærboelling I, Vesth TC, Frisvad JC, Nybo JL, Theobald S, Kuo A, et al. Linking secondary
565 metabolites to gene clusters through genome sequencing of six diverse *Aspergillus* species.
566 *Proc Natl Acad Sci.* 2018;115:E753–61.

- 567 29. Quick J. Ultra-long read sequencing protocol for RAD004 v3.
568 <https://www.protocols.io/view/ultra-long-read-sequencing-protocol-for-rad004-mrxc57n>.
569 Accessed 25 Oct 2019
- 570 30. Wick R. Porechop. <https://github.com/rrwick/Porechop>. Accessed 3 Dec 2019
- 571 31. Marijon P, Chikhi R, Varré J-S. yacrd and fpa: upstream tools for long-read genome
572 assembly. *bioRxiv*. 2019;674036.
- 573 32. Kielbasa SM, Wan R, Sato K, Horton P, Frith MC. Adaptive seeds tame genomic sequence
574 comparison. *Genome Res*. 2011;21:487–93.
- 575 33. Krzywinski M, Schein J, Birol Ī, Connors J, Gascoyne R, Horsman D, et al. Circos: An
576 information aesthetic for comparative genomics. *Genome Res*. 2009;19:1639–45.
- 577 34. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.
578 2018;34:3094–100.
- 579 35. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
580 Alignment/Map format and SAMtools. *Bioinforma Oxf Engl*. 2009;25:2078–9.
- 581 36. Thorvaldsdóttir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-
582 performance genomics data visualization and exploration. *Brief Bioinform*. 2013;14:178–92.
- 583 37. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and
584 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*.
585 2017;27:722–36.
- 586 38. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
587 *J Mol Biol*. 1990;215:403–10.
- 588 39. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
589 improvements in performance and usability. *Mol Biol Evol*. 2013;30:772–80.
- 590 40. Gouy M, Guindon S, Gascuel O. SeaView Version 4: A Multiplatform Graphical User
591 Interface for Sequence Alignment and Phylogenetic Tree Building. *Mol Biol Evol*.
592 2010;27:221–4.
- 593 41. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly from
594 long uncorrected reads. *Genome Res*. 2017;27:737–46.
- 595 42. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile
596 and open software for comparing large genomes. *Genome Biol*. 2004;5:R12.
- 597 43. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
598 data. *Bioinformatics*. Oxford Academic; 2014;30:2114–20.
- 599 44. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
600 *ArXiv* 2013;1303.3997

601 45. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An
602 Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly
603 Improvement. PLoS One. 2014;9:e112963.

604

605

606

Figure legends

Figure 1. An overview of *D. coniospora* strain isolation and culture history.

A strain of *D. coniospora* collected from Denmark in 1982 at the latest was deposited at the CBS-KNAW culture collection, now held by the Westerdijk Fungal Biodiversity Institute as CBS615.82. It was transferred in 1987 to the ARS Collection of Entomopathogenic Fungal Cultures (as ARSEF 2468) and then re-isolated in 2001 as ARSEF 6962. A second strain collected from Sweden was deposited at the American Type Culture Collection as ATCC 96282. It has been cultured through serial passage in *C. elegans* continuously since 1999.

Figure 2. Inter-chromosomal rearrangements between strains Swe2 and Dan2.

A. Circos plot representing regions >6 kb that are very similar between Dan2 (left – olive) and Swe2 (right – light blue) assemblies as determined by an all-against-all LAST analysis. Swe2 contig numbers are the last two digits of the accession ID (shown in B), preceding the suffix. Red and dark blue rectangles represent rearrangement junctions probed by PCR. B. Conceptual design of the PCR primers. C. Amplicons from the PCR were visualized after electrophoresis. Each pair gave one specific band of the expected size. The colour code is the same for the 3 panels.

Figure 3. Comparisons between Canu and Flye assemblies.

A, B. Dot plots of the non-congruent assemblies generated by Canu (x-axis) against those generated by Flye (y-axis) for the Swe3 (A) and Dan1 (B) genomes. The orange triangle (A) marks the position where the Canu contig *tig00000004* was split during the manual curation because of its chimeric nature. The green arrow (B) marks the position

of a Flye scaffolding error. C. Schematic representation of the Dan1 Flye assembly, showing the mapping of chimeric reads close to the scaffolding error (green triangle). The coordinates in brackets are the mapping positions of the clipped part of the reads (dash line) on another contig of the assembly. Notably, this error was eliminated when these chimeric reads were excluded from the input data. D. Mapping of long-reads close to the scaffolding error (green triangle) on the Dan1 Flye assembly. The green bar marks the telomeric tandem repeat motif. The grey bar indicates the 100 Ns inserted by Flye to unite the scaffold.

Figure 4. Synteny among the genomes of 5 *D. coniospora* strains.

Circos plot representing regions >20 kb that are very similar between assemblies as determined by all-against-all LAST analyses. Each assembly is shown at the same scale and in the same order and orientation across panels.

Figure 5. Evaluation of sequence errors in the 3 new genomes.

A. Percentage of correct genes (based on length of the corresponding predicted protein) among 305 conserved genes, for the 3 new genomes, in the initial assembly and after two different polishing strategies. B, C, D. Scatter plots of homopolymer composition (A/T or C/G) and accuracy among the same 305 conserved genes for Dan1 (B), Swe1 (C) and Swe3 (D). The dot size is proportional to the number of genes, and the colour indicates the proportion of genes predicted to be correct. Red and purple arrows highlight two particular cases, among many, where homopolymer errors are only present in one genome.

Figure 6. Sequence anomalies introduced by assembly and/or polishing tools.

A. A comparison of one small region of the Swe3 sequence before (top) and after polishing (Racon x4 and Medaka; bottom). As indicated by the orange line, long stretches of A and T homopolymers are introduced by polishing, in the absence of coherent read support. B. From top to bottom, the assembly produced by Canu excludes a region of around 10 kb, despite strong read support. After 2 and 4 iterations, Racon progressively filled the gap. Medaka then introduced an insert of roughly the correct size, but of aberrant sequence composition. For each panel, the height of the boxes in the top line indicates the read coverage for each base. A grey box indicates full agreement with the consensus sequence, otherwise the colour indicates the proportion of read support for each nucleotide (G, tan; C, blue; A, green; T, red). Below this, the ONT reads that align in forward (pink) and reverse (blue) orientation are shown as lines. A coloured letter or purple rectangle show a difference (nucleotide variant or insertion in reads, respectively) in the read's sequence compared to the genome sequence. C. The 10 kb sequence introduced by polishing is of aberrant composition as illustrated by the region immediately surrounding the 5' breakpoint (yellow arrowhead). There are single nucleotide errors introduced despite coherent read support for the "Before" sequence (light blue dots), and then a continuous stretch, exemplified by A and T homopolymers that lack any sequence support at all (light blue line).

Figure 7. Example of sequence errors introduced during assembly and polishing.

A. Stacked bar plot of USCO status for the orthologues of selected mono-exonic Swe2 or Dan2 genes classified according to the result of a TBLASTN search against the

indicated assembly (Canu: from Canu; Racon: after long-read polish; Pilon: after short-read polishing). B. Detailed view of 2 parts of RJ55_06485 from the Swe2 reference genome each containing a homopolymer sequence (underlined) and the corresponding positions in successive Swe1 assemblies. For each, the predicted protein sequence, highlighted in turquoise, with the other open reading frames in grey, is shown above the corresponding nucleotide sequence. The red arrow heads highlight the missing nucleotides, the extraneous nucleotide is boxed in red.

Figure 1

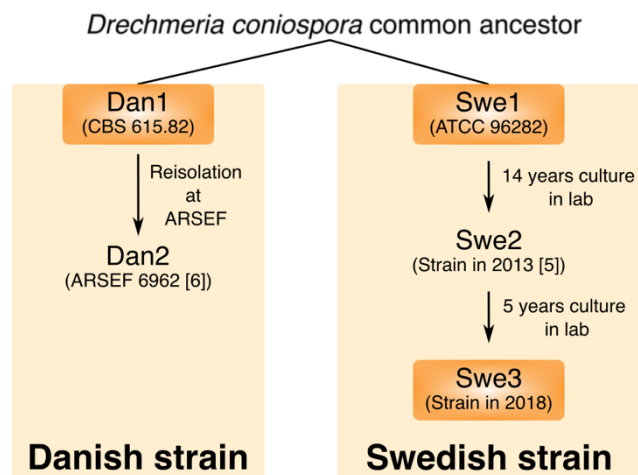


Figure 2

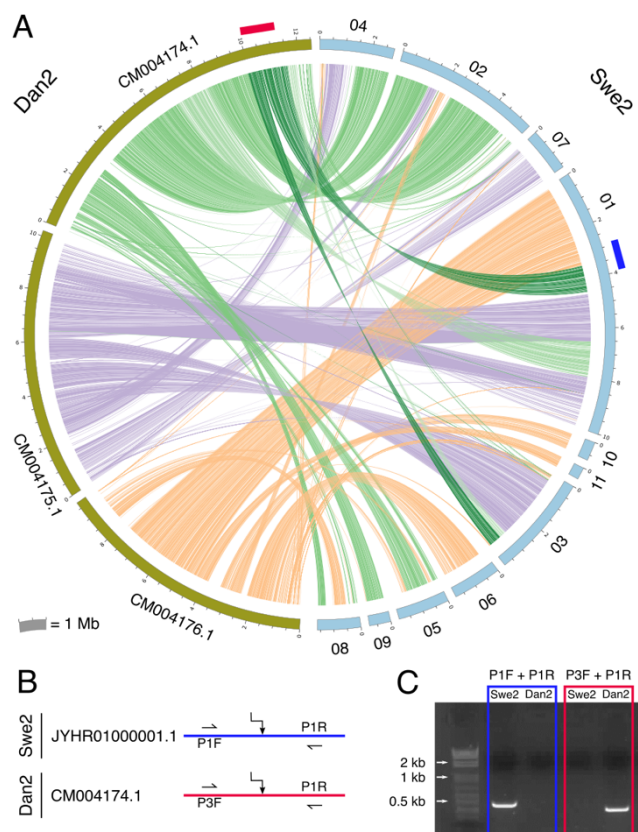


Figure 3

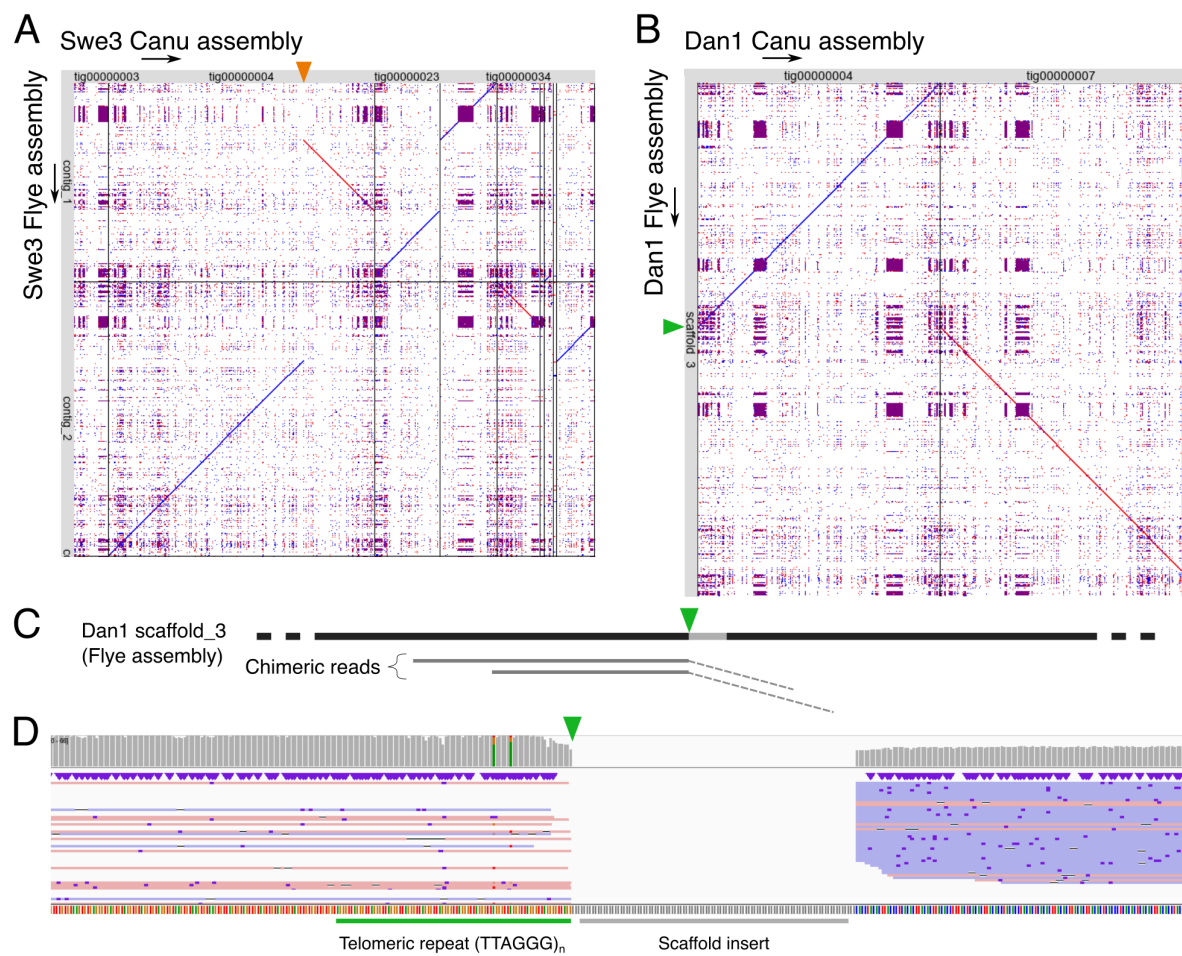


Figure 4

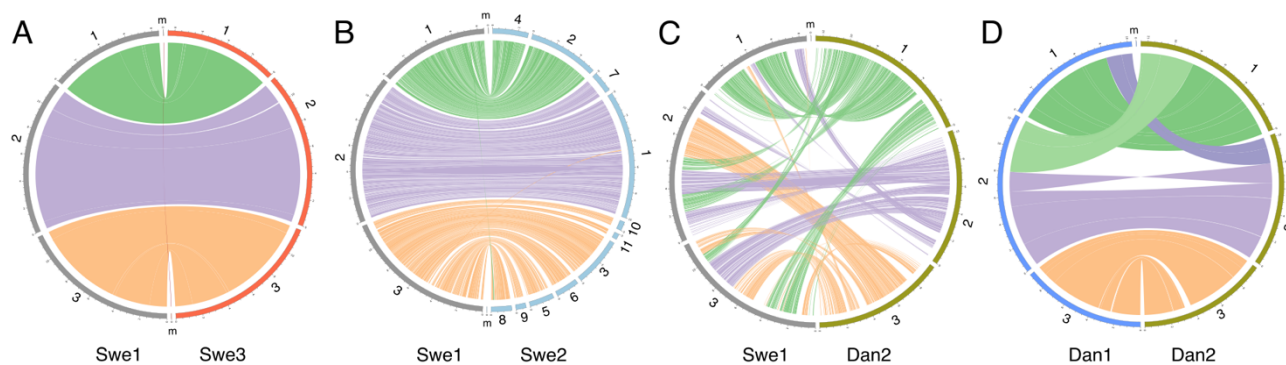


Figure 5

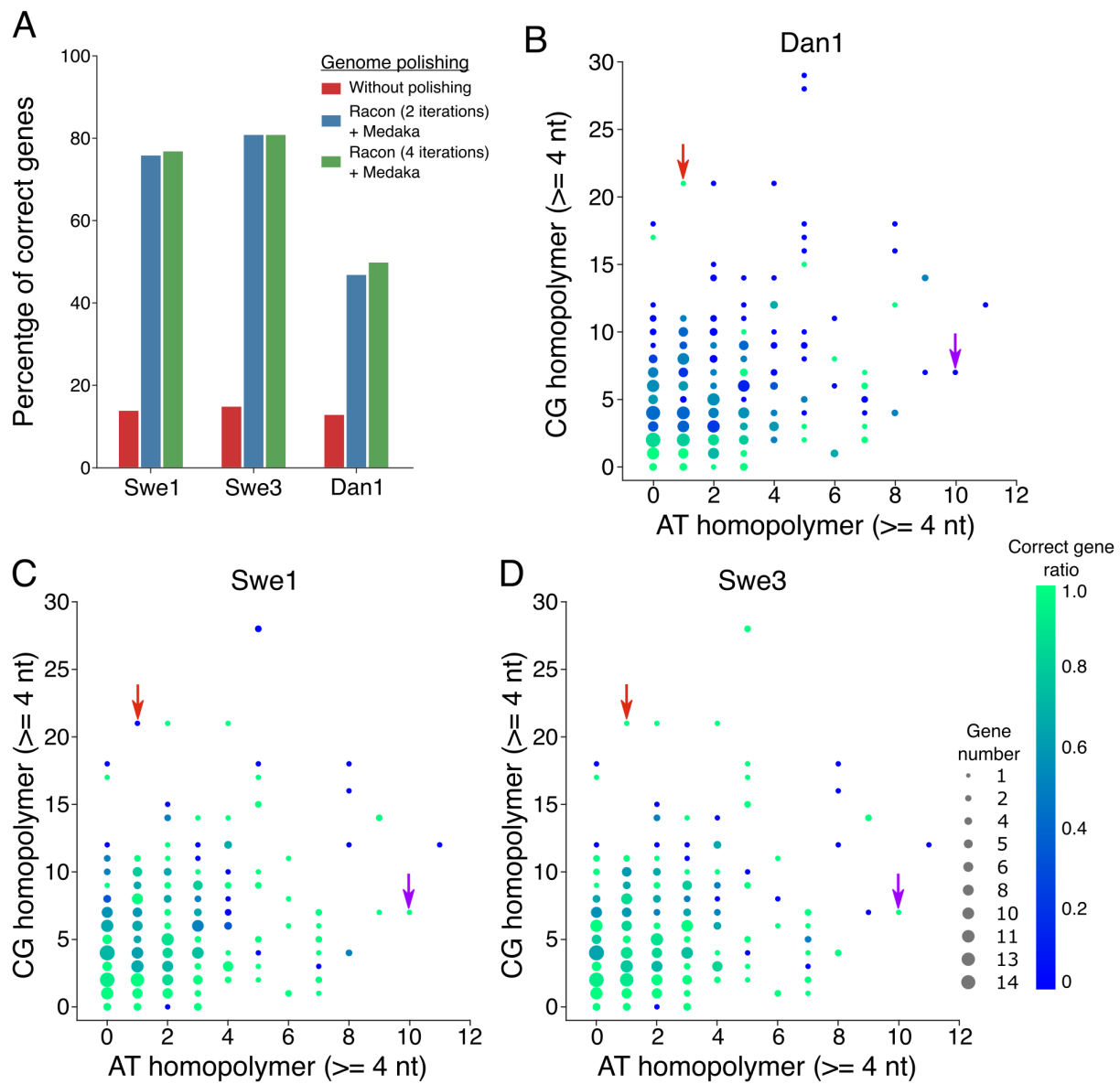


Figure 6

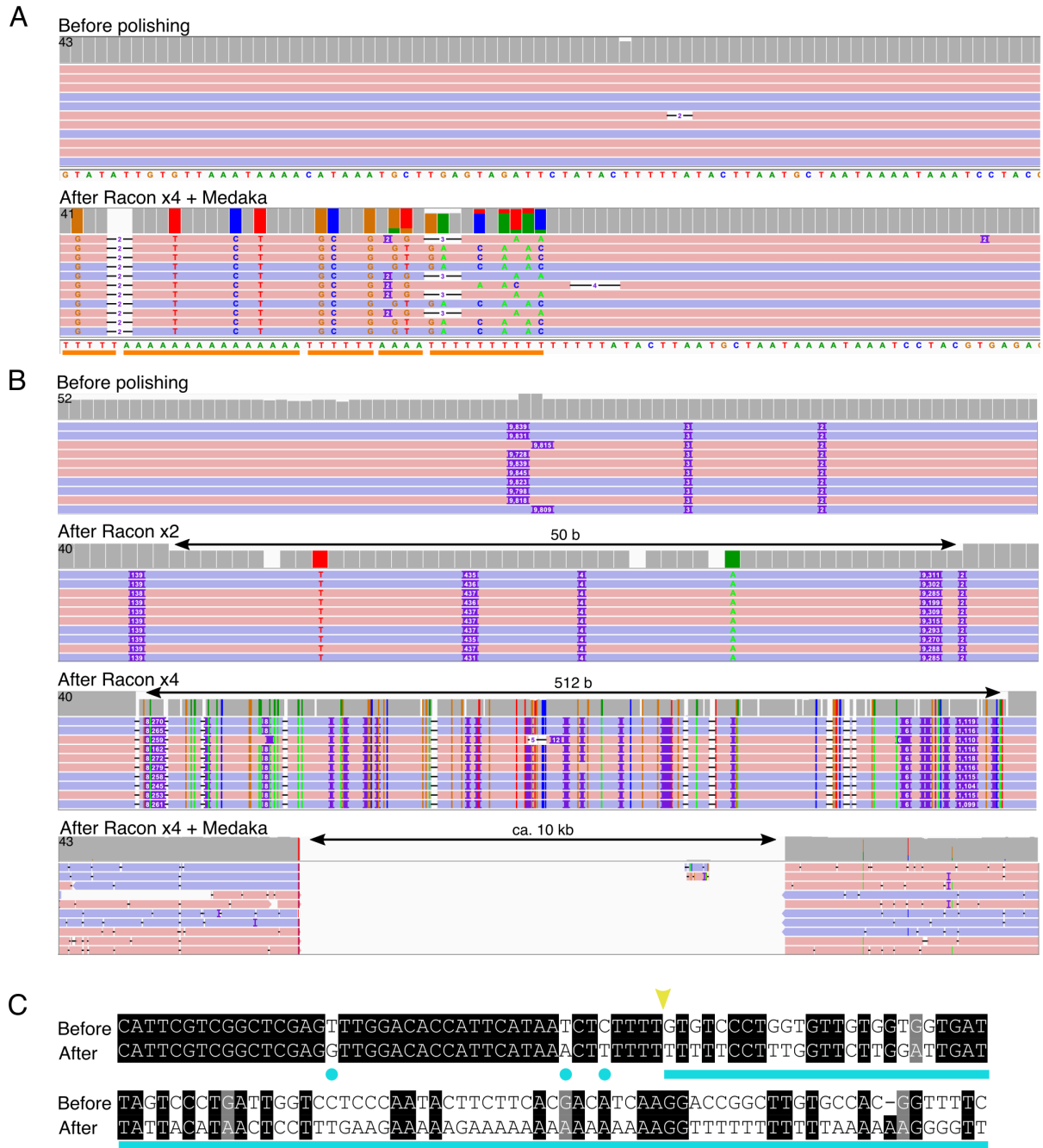
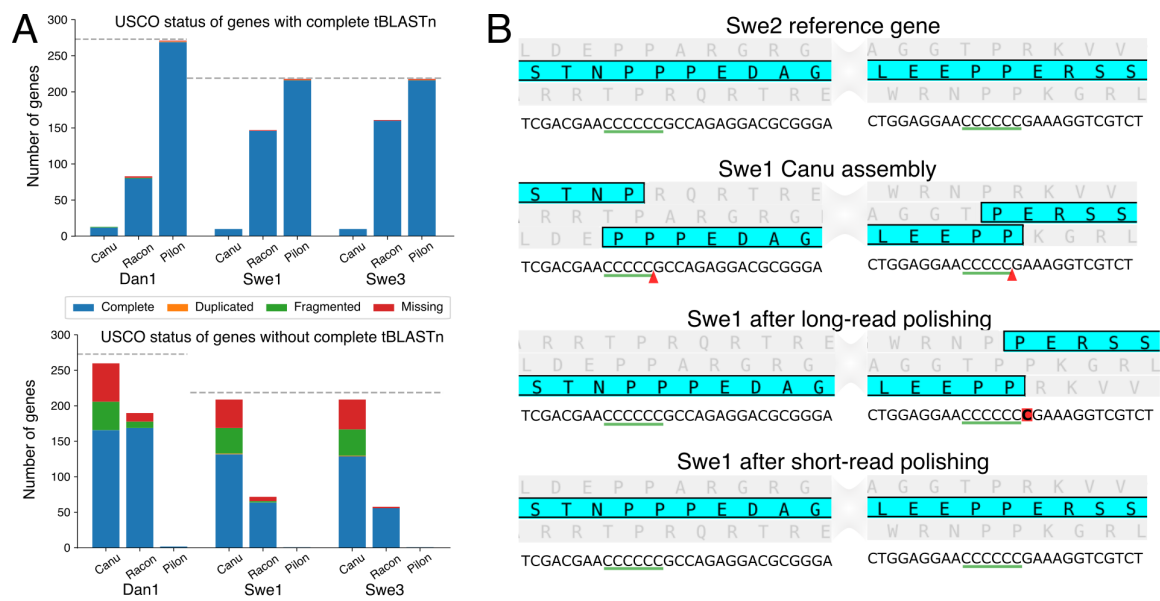
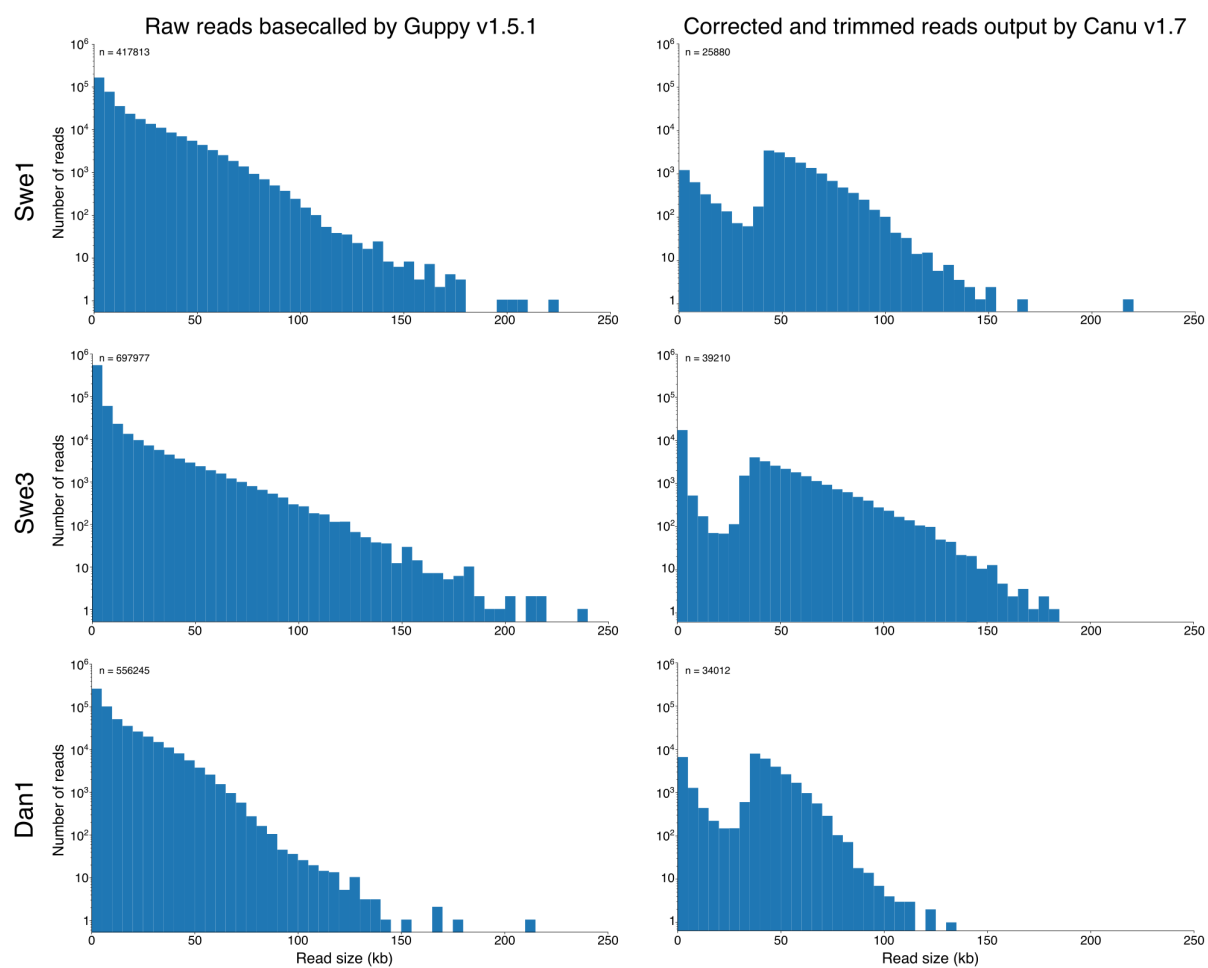


Figure 7



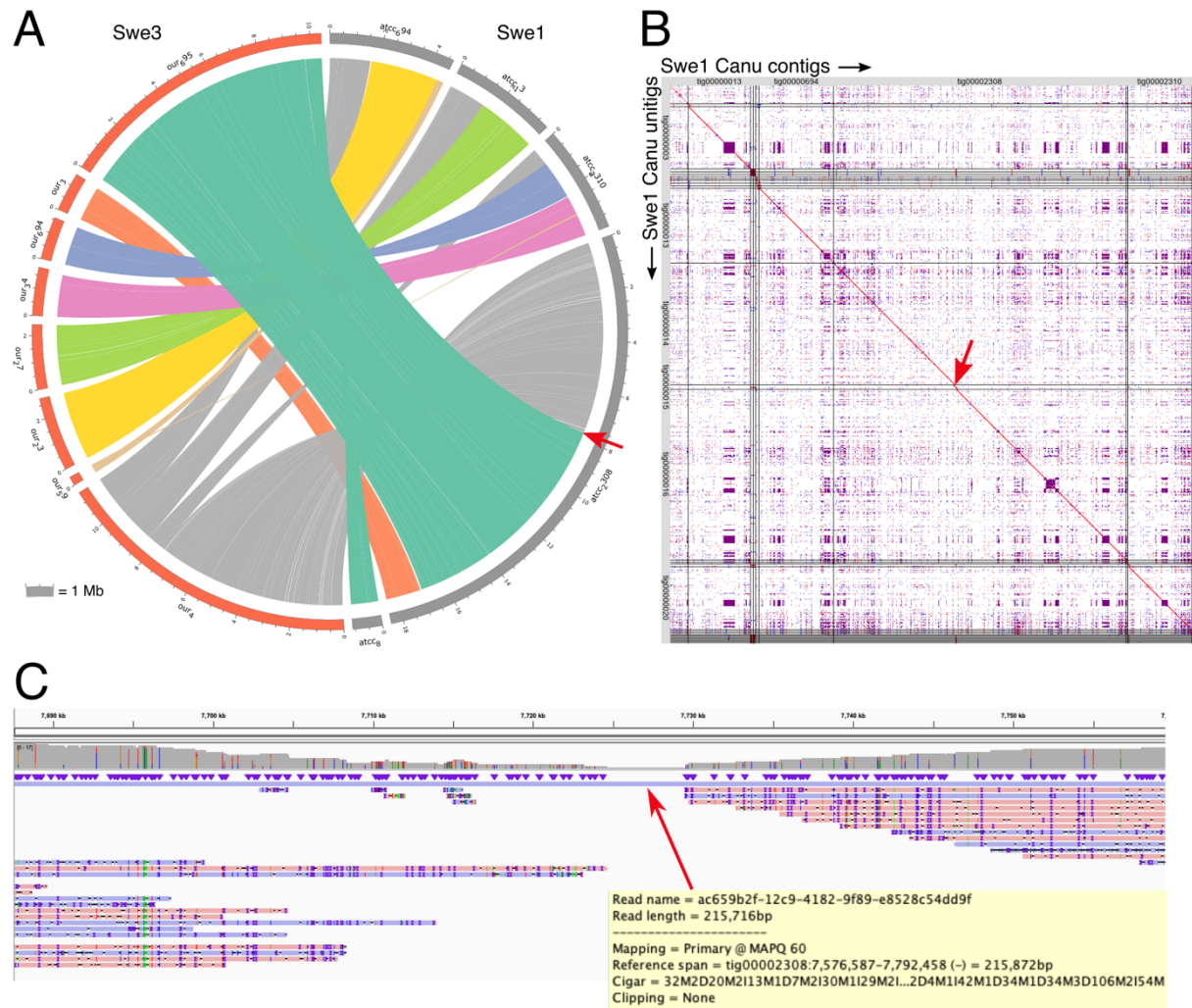
Supplementary Figure 1



Supplementary Figure 1. Distribution of the size of the reads and the assembly statistics.

Distributions in 5 kb bins of the size of the set of reads basecalled by Guppy v1.5.1 (left-hand panels), and of reads corrected and trimmed by Canu for the initial assemblies (righthand panels).

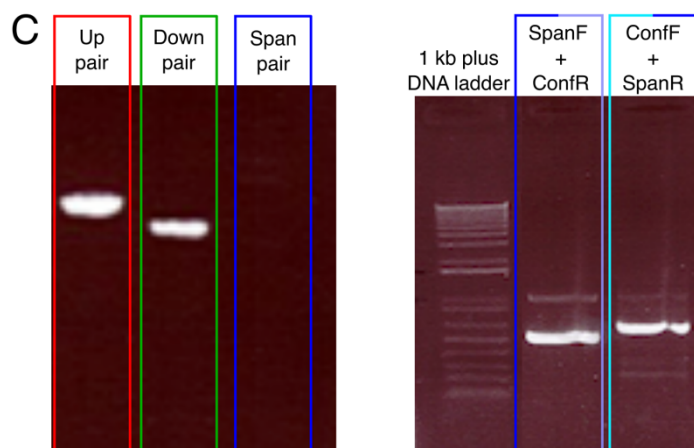
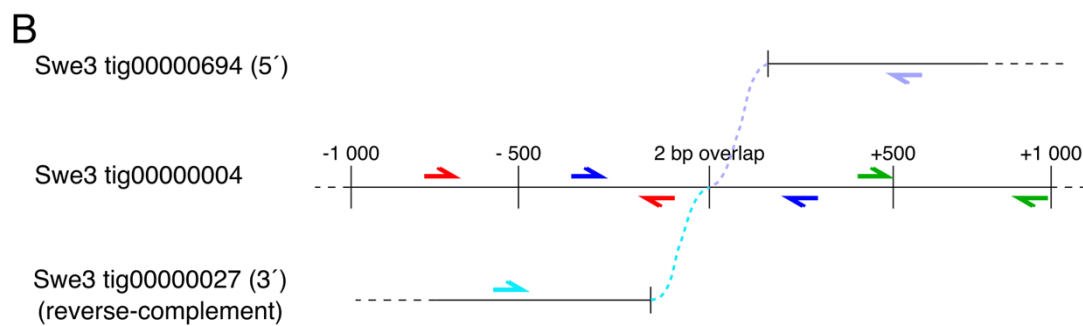
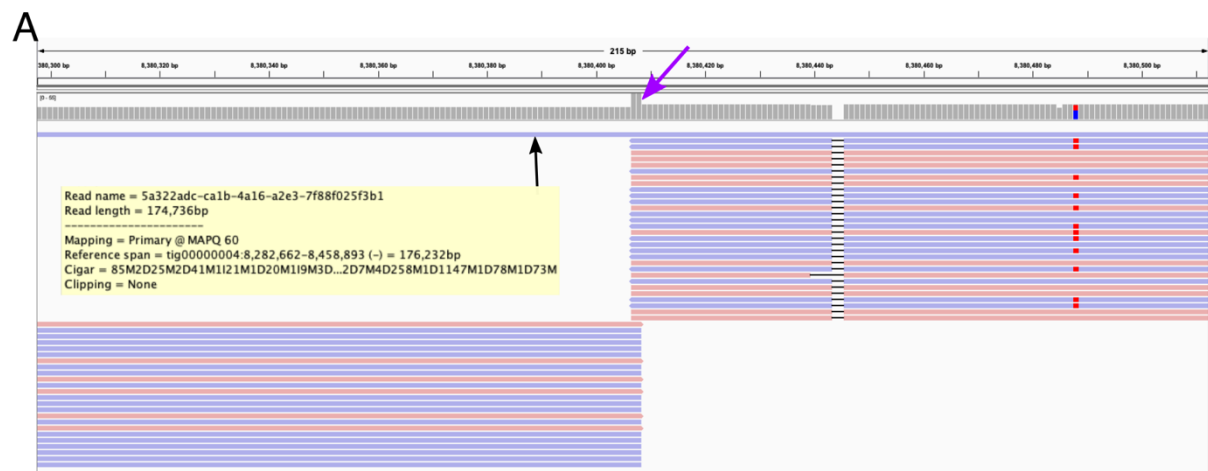
Supplementary Figure 2



Supplementary Figure 2. Detection of an error in the initial Swe1 assembly introduced by a long chimeric read.

A. Circos plot representing regions >20 kb that are very similar between Swe3 and Swe1 Canu assemblies as determined by an all-against-all LAST analysis. The red arrow indicates a break in the synteny for the largest Swe1 contig. B. Dot-plot of an all-against-all comparison of the Swe1 contigs and unitigs produced by Canu. Contigs are contiguous sequences present in the primary assembly, including both unique and repetitive elements. Unitigs are contigs split at alternate paths in the assembly graph. The red arrow indicates the discontinuity in the alignment between Swe1 contigs and unitigs. This occurs on the same contigs and at the same coordinates as in A. C. Mapping of the Swe1 reads, corrected and trimmed by Canu, on the Swe1 Canu assembly (detail of around 70 kb on contig tig00002308 flanking the synteny break). Moving into the central 5 kb region, read support progressively drops from 40 to just one, corresponding to a very long read of about 215 kb, which was shown to be chimeric.

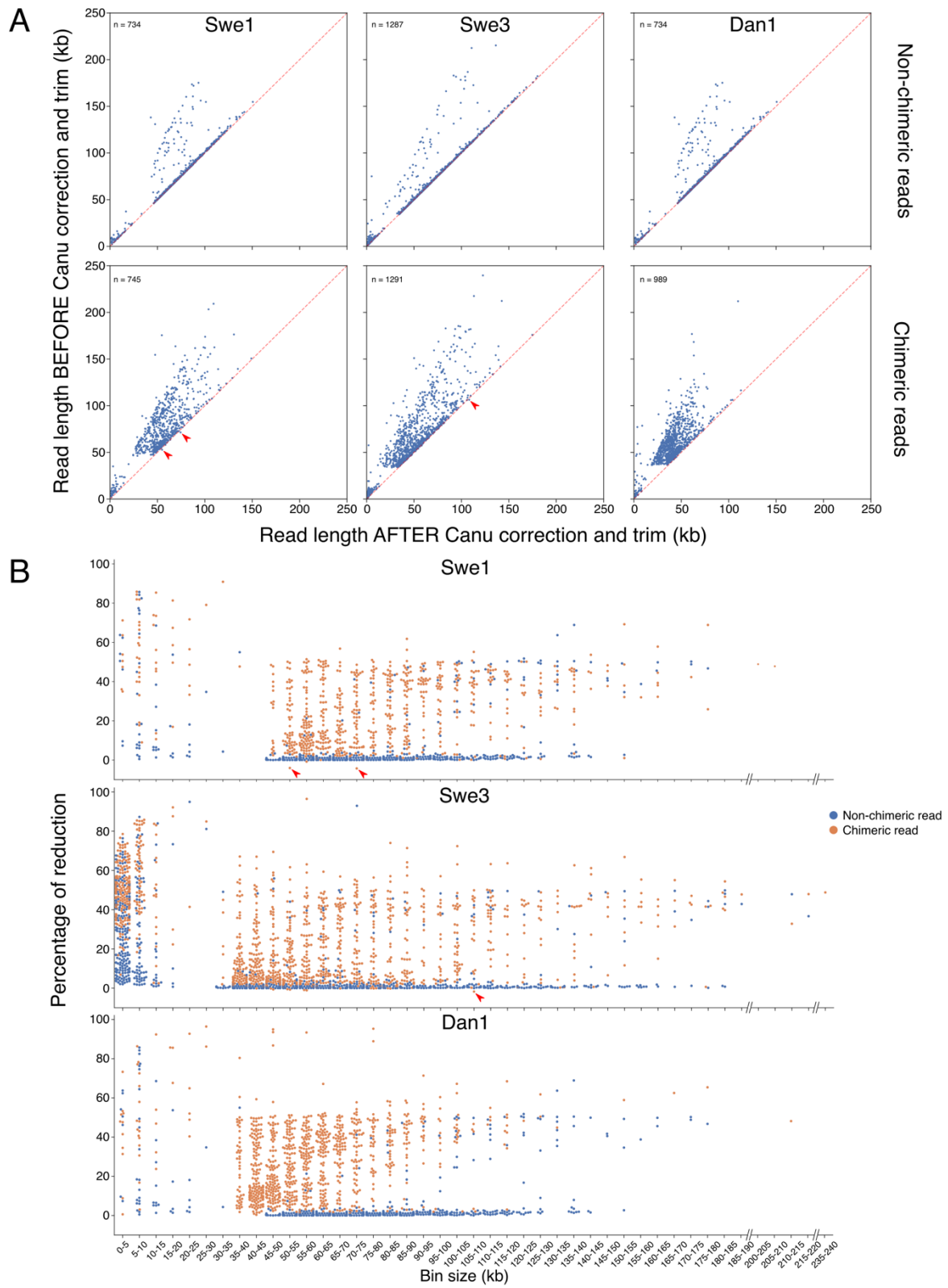
Supplementary Figure 3



Supplementary Figure 3. Detection of an error in the Swe3 initial assembly introduced by a long chimeric read.

A. Mapping of the Swe3 reads, corrected and trimmed by Canu, on the Swe3 Canu assembly. The putative 175 kb chimeric read was identified by a break of synteny (not shown) and because of a sharp (ca. 2-fold) increase in coverage, spanning only 2 nucleotides, (purple arrow) on the contig tig0000004. B. Conceptual design of the PCR primers used to verify the assembly. Three pairs of primers were designed on the tig0000004: the Up pair (red), the Down pair (green) and the Span pair (blue). Two other primers were designed on the basis of the corrected assembly sequence (dotted coloured lines): SpanF in turquoise on the contig tig00000027 and SpanR in light purple on the contig tig00000695. C. PCR result of the different pairs used. Amplicons had the expected sizes.

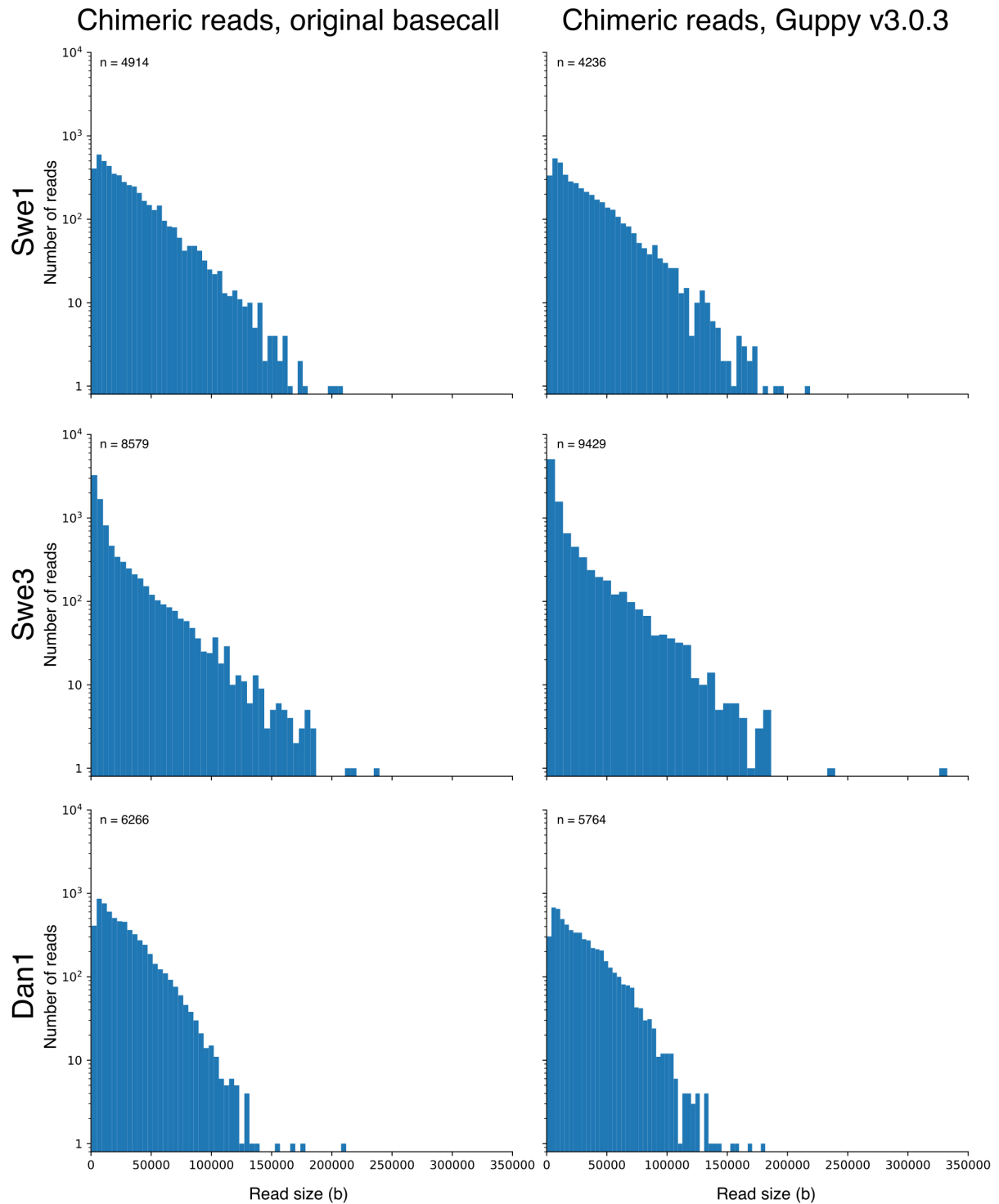
Supplementary Figure 4



Supplementary Figure 4. Canu correction of non-chimeric and chimeric reads.

A. Scatter plots of the size of reads in the initial dataset (y-axis) and after Canu correction (x-axis) for non-chimeric (top) and chimeric (bottom) reads. Subsets of non-chimeric reads of a similar number (as indicated) and having the equivalent size distribution (before Canu correction, +/- 1% in length) as the chimeric read were used to allow a valid comparison. B. Percentage of read length reduction after Canu correction, in 5 kb bins. The bin size corresponds to the read size in the raw dataset. The red arrows highlight the rare chimeric reads that are longer after the correction step.

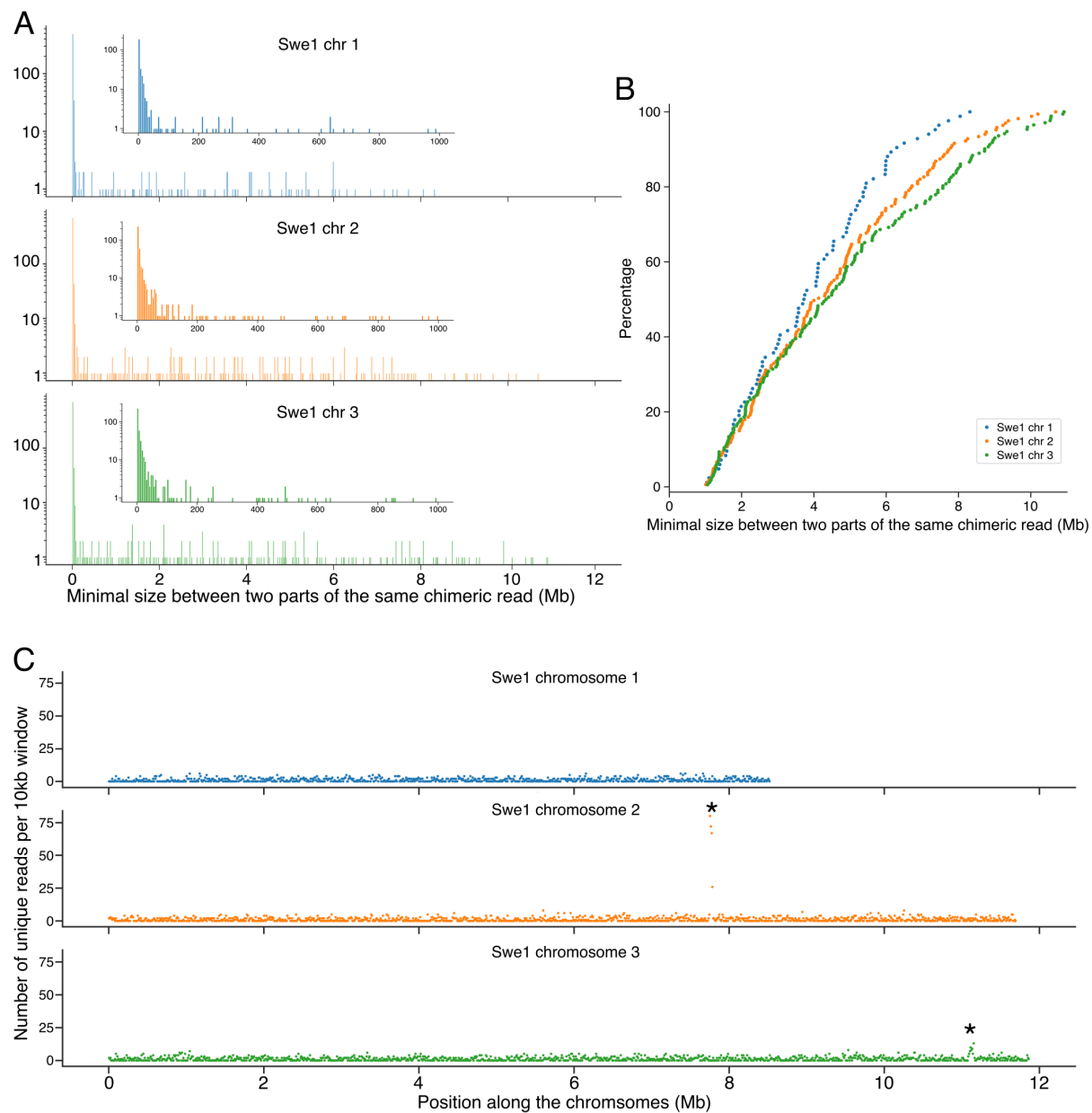
Supplementary Figure 5



Supplementary Figure 5. Size distribution of chimeric reads.

Distributions in 5 kb bins of the size of the set of reads identified as chimeric by YACRD (see Methods), in the original dataset (left panels) and amongst reads basecalled by Guppy v3.0.3 (right panels).

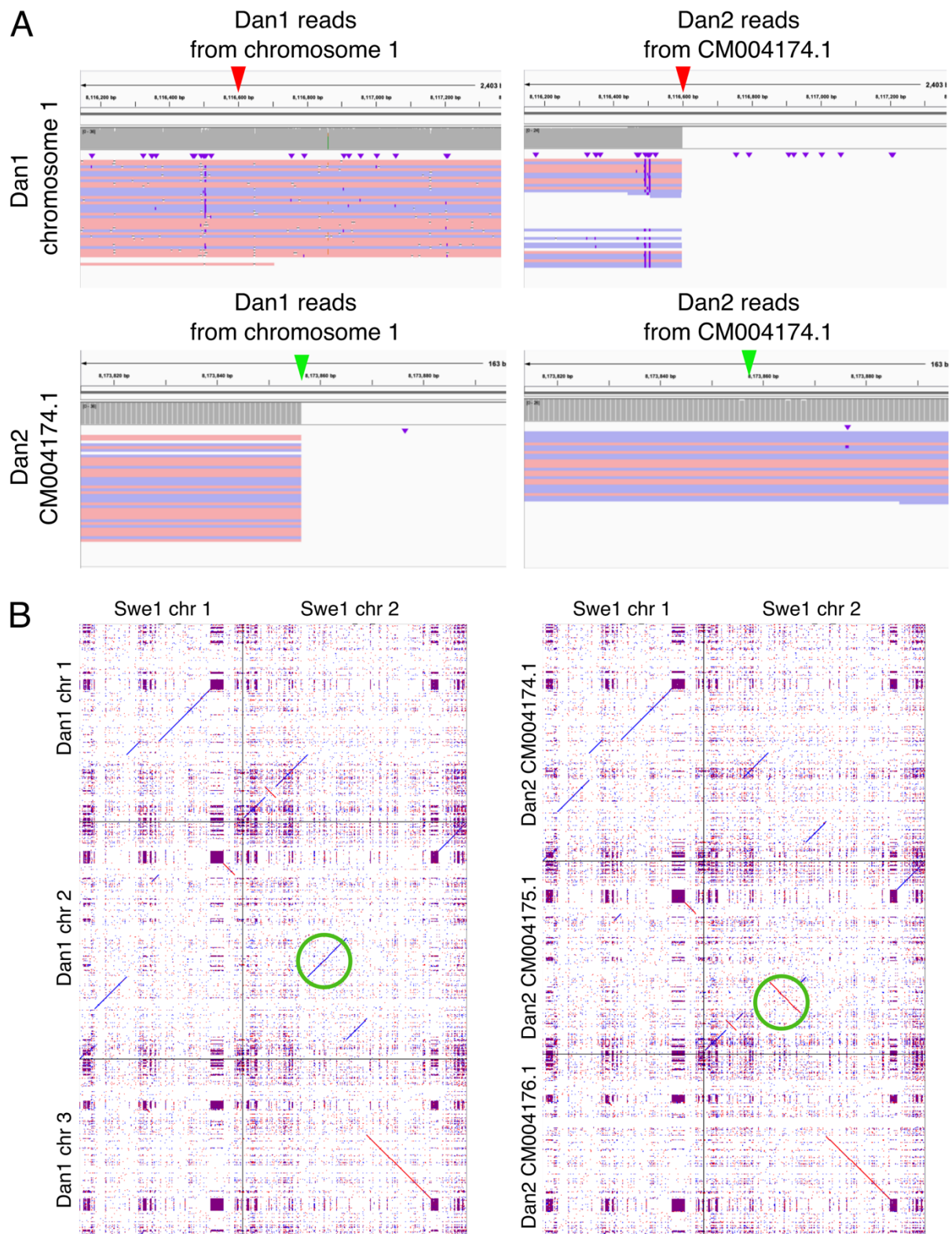
Supplementary Figure 6



Supplementary Figure 6. Analysis of the Swe1 chimeric reads.

For each chimeric read, the 500 bp of sequence on each side of the putative break point was mapped onto the final Swe1 genome. A, Distribution in 500 equal bins of the interval between the two parts of each chimeric read. The inserts highlight reads that map within 1 kb of each other (in 200 bins). B, Cumulative distribution of the interval between the two parts of each chimeric read (for those separated by > 1 Mb). The distance separating the two parts is broadly spread. C. Distribution of mapped sequences from each side of the putative breakpoint for each chimeric read that mapped to a single chromosome. The peak on chromosome 2 corresponds to the site where the nuclear genome matches the mitochondrial DNA, and that on chromosome 3 to a highly repeated region containing tandem copies of rDNA. In both cases, these therefore reflect erroneous attribution of chimerism. There are thus no true hot spots for chimeric read breakpoints on any chromosome.

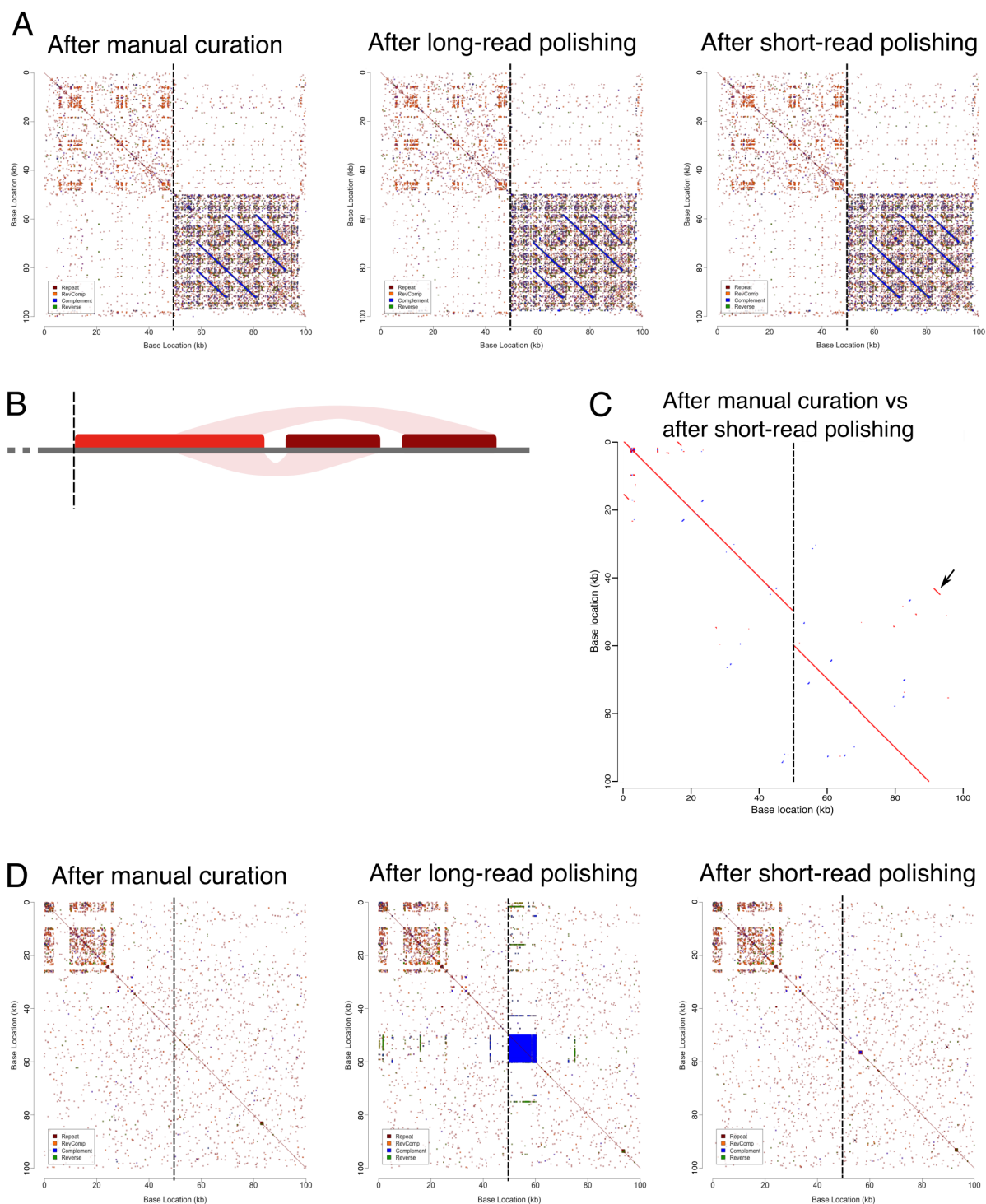
Supplementary Figure 7



Supplementary Figure 7. Identification of intra- and inter-chromosomal rearrangements between Dan1 and Dan2.

A. Mapping of long reads from Dan1 (left panels) and Dan2 (right panels) on Dan1 chromosome 1 (top panels) and Dan2 chromosome 1 (CM004174.1; bottom panels). The arrowheads highlight points of discontinuity in the read coverage, consistent with chromosomal rearrangements between Dan1 and Dan2. B. Alignment of selected chromosomes of the final assemblies between Dan1 and Swe1 (left) and Dan2 and Swe1 (right). The unique difference in the orientation of part of Swe1 chromosome 2 between Dan1 and Dan2 is highlighted by the green circles.

Supplementary Figure 8



Supplementary Figure 8. Analysis of the pattern of repeated sequences in the neighbourhood of polishing error in the Swe3 genome.

A. Dot-plot based on k-mer identity ($k = 11$; See Supplementary Methods) of the neighbourhood of the *numts* vs itself, from the left to the right, before any polishing, after long-read polishing and after short-read polishing. The vertical dashed line at the centre of each plot indicates the start of the *numts*. B. A schematic representation of the large-scale genomic organization in the same region. The rectangles (red and dark red) represent *numts* (identified by BLASTN against the Swe3 mitochondrial genome) and the ribbons show the regions of sequence similarity (>99.5%) between them. For the Swe1 assembly, one dark-red rectangle is absent, while in the Swe2 assembly, the two dark-red rectangles are absent. C. An alignment of the area before any polishing (x-axis) and after short-read polishing (y-axis) reveals the re-inclusion of this missing sequence in the final assembly. The arrow highlights a 1.2 kb region that is duplicated on both sides of the sequence discontinuity (vertical dashed line). D. Dot-plots based on k-mers ($k = 11$) of the same region. The blue square reflects the presence of a region of low-complexity repeated sequences. The pattern at the top left hand of the plots reflect a region of low-complexity sequence that is not shared across the sequence discontinuity (vertical dashed line).

Supplementary table 1: Coverage, genome length and results of the tBLASTn

	Set of reads	Swe1	Swe3	Dan1
Coverage (fold)	Long read - Guppy 1.5.1	136	94	150
	Long read - Canu corrected	36	36	36
	Long read - Guppy 3.0.3	140	97	155
	Short read	406	430	482
Genome length (bp)	Canu after curation	31971578	32076624	31919747
	Long-read polishing	32121045	32271010	32120031
	Short-read polishing	32117370	32271155	32095731
TBLASTN result for the 305 benchmarking genes	Canu after curation	43 (14%)	48 (16%)	42 (14%)
	Long-read polishing	236 (77%)	248 (81%)	152 (50%)
	Short-read polishing	305 (100%)	305 (100%)	305 (100%)

The coverage for each set of reads was computed as the mean depth for the assemblies after short-read polishing, using Samtools coverage v1.10.

Supplementary Table 2: Examples of read support for different homopolymer sequences in 10 representative genes.

gene	reference	Swe2	Swe1	Swe1 read	Swe1 guppy remark	Swe3	Swe3 read	Swe3 guppy remark	Dan2	Dan1	Dan1 read	Dan1 guppy remark
ODA80475.1	c	5	5	36 (2)	supported	5	34 (0)	supported	5	6	28 (0)	weak support (<10% read)
ODA80475.1	t	4	4	35 (0)	supported	4	34 (3)	supported	4	4	28 (12)	supported
ODA80475.1	t	NA	NA	NA (NA)	NA	NA	NA (NA)	NA	4	4	28 (3)	supported
ODA80475.1	c	4	4	37 (23)	supported	4	35 (33)	supported	4	5	27 (2)	supported
ODA80475.1	g	7	8	36 (0)	weak support (<10% read)	6	35 (2)	supported	7	7	28 (0)	supported
ODA80133.1	c	4	4	33 (33)	supported	4	41 (41)	supported	4	4	37 (36)	supported
ODA80133.1	c	4	4	33 (26)	supported	4	41 (38)	supported	4	4	36 (18)	supported
ODA80133.1	g	4	4	33 (32)	supported	4	42 (40)	supported	4	4	35 (35)	supported
ODA80133.1	c	7	6	32 (0)	weak support (<10% read)	6	41 (0)	supported	7	7	37 (0)	supported
ODA83792.1	c	4	4	33 (33)	supported	4	38 (38)	supported	4	4	28 (28)	supported
ODA83792.1	c	4	4	33 (30)	supported	4	38 (36)	supported	4	4	28 (24)	supported
ODA83792.1	c	4	4	33 (32)	supported	4	38 (34)	supported	4	4	28 (28)	supported
ODA83792.1	c	5	6	33 (0)	weak support (<10% read)	5	38 (17)	supported	5	6	28 (0)	weak support (<10% read)
ODA83792.1	c	4	4	33 (33)	supported	4	38 (36)	supported	4	4	28 (27)	supported
ODA83792.1	c	4	4	32 (32)	supported	4	38 (36)	supported	4	4	28 (28)	supported
ODA83792.1	c	5	5	32 (0)	supported	5	38 (0)	supported	5	5	28 (1)	supported
ODA83792.1	c	5	5	32 (12)	supported	6	38 (0)	weak support (<10% read)	5	5	28 (11)	supported
ODA83792.1	c	4	4	32 (32)	supported	4	38 (37)	supported	4	4	28 (28)	supported
ODA83792.1	c	4	4	32 (31)	supported	4	38 (36)	supported	4	4	28 (28)	supported
ODA83792.1	a	4	4	33 (28)	supported	4	37 (35)	supported	4	4	28 (27)	supported
ODA83792.1	a	4	4	33 (28)	supported	4	37 (37)	supported	4	4	28 (21)	supported
ODA83792.1	g	4	4	32 (32)	supported	4	38 (37)	supported	4	4	28 (28)	supported
ODA83792.1	t	4	4	32 (32)	supported	4	38 (38)	supported	4	4	28 (28)	supported
ODA83726.1	c	4	5	25 (0)	weak support (<10% read)	4	37 (12)	supported	4	5	37 (0)	supported
ODA83726.1	c	4	4	25 (23)	supported	4	37 (35)	supported	4	4	37 (36)	supported
ODA83726.1	c	5	5	25 (17)	supported	5	37 (27)	supported	5	5	37 (14)	supported
ODA83726.1	c	5	6	25 (0)	weak support (<10% read)	5	37 (12)	supported	5	6	37 (0)	supported
ODA83726.1	c	7	7	25 (1)	weak support (<10% read)	7	37 (0)	weak support (<10% read)	7	7	37 (2)	supported
ODA83726.1	c	5	5	25 (24)	supported	5	36 (28)	supported	5	5	37 (26)	supported
ODA83726.1	c	5	5	25 (6)	supported	5	36 (13)	supported	5	5	37 (14)	supported
ODA83726.1	g	4	4	25 (24)	supported	4	36 (29)	supported	4	4	37 (20)	supported
ODA84322.1	c	5	5	38 (1)	supported	5	41 (0)	supported	5	5	41 (11)	supported
ODA84322.1	g	4	4	38 (29)	supported	4	41 (29)	supported	4	4	41 (21)	supported
ODA84322.1	g	4	4	38 (13)	supported	4	40 (7)	supported	4	5	41 (1)	supported
ODA84322.1	c	5	5	38 (6)	supported	5	42 (3)	supported	5	5	41 (3)	supported
ODA84322.1	c	5	5	38 (6)	supported	5	42 (6)	supported	5	5	41 (1)	supported
ODA84322.1	c	5	6	38 (0)	weak support (<10% read)	6	42 (3)	weak support (<10% read)	5	5	41 (7)	supported
ODA84322.1	g	4	4	37 (37)	supported	4	42 (42)	supported	4	5	41 (1)	supported
ODA82449.1	c	5	5	39 (4)	supported	5	46 (2)	supported	5	5	26 (8)	supported
ODA82449.1	c	5	6	39 (1)	weak support (<10% read)	6	46 (1)	weak support (<10% read)	5	8	26 (0)	weak support (<10% read)
ODA82449.1	c	5	5	39 (16)	supported	6	46 (2)	weak support (<10% read)	5	5	26 (8)	supported
ODA82449.1	c	5	6	39 (0)	weak support (<10% read)	6	46 (0)	weak support (<10% read)	5	6	26 (0)	weak support (<10% read)
ODA82449.1	a	5	5	38 (21)	supported	5	47 (34)	supported	5	5	26 (6)	supported
ODA82449.1	c	6	7	38 (0)	weak support (<10% read)	7	47 (2)	weak support (<10% read)	6	6	26 (2)	supported
ODA82449.1	c	5	5	38 (25)	supported	5	47 (33)	supported	5	5	26 (15)	supported
ODA81732.1	g	4	5	19 (0)	supported	5	40 (0)	supported	4	5	39 (0)	supported
ODA81732.1	c	7	7	20 (1)	weak support (<10% read)	7	40 (2)	weak support (<10% read)	7	8	40 (1)	weak support (<10% read)
ODA76966.1	c	5	5	40 (1)	supported	5	31 (1)	supported	5	6	39 (2)	weak support (<10% read)
ODA76966.1	g	4	4	40 (40)	supported	4	31 (30)	supported	4	4	39 (39)	supported
ODA76966.1	c	4	4	40 (32)	supported	4	31 (30)	supported	4	4	39 (29)	supported

ODA76966.1	g	4	4	4	39 (37)	supported	4	31 (30)	supported	4	4	39 (37)	supported
ODA76966.1	c	4	4	4	39 (31)	supported	4	31 (28)	supported	4	4	39 (29)	supported
ODA76966.1	g	8	7	38 (0)	weak support (<10% read)	weak support (<10% read)	7	31 (1)	weak support (<10% read)	8	7	39 (0)	weak support (<10% read)
ODA76966.1	g	4	4	38 (38)	supported	supported	4	31 (31)	supported	4	4	39 (39)	supported
ODA76653.1	c	4	4	39 (37)	supported	supported	4	38 (36)	supported	4	4	40 (39)	supported
ODA76653.1	a	4	4	39 (39)	supported	supported	4	38 (38)	supported	4	4	40 (40)	supported
ODA76653.1	g	4	4	39 (39)	supported	supported	4	38 (37)	supported	4	4	40 (40)	supported
ODA76653.1	a	4	4	39 (39)	supported	supported	4	38 (38)	supported	4	4	40 (40)	supported
ODA76653.1	g	4	4	39 (37)	supported	supported	4	38 (8)	supported	4	4	40 (38)	supported
ODA78969.1	c	4	4	24 (22)	supported	supported	4	29 (28)	supported	4	4	42 (38)	supported
ODA78969.1	c	4	4	23 (23)	supported	supported	4	29 (29)	supported	4	4	42 (42)	supported
ODA78969.1	c	4	4	23 (17)	supported	supported	4	29 (20)	supported	4	4	42 (19)	supported
ODA78969.1	c	4	4	23 (22)	supported	supported	4	29 (26)	supported	4	5	42 (2)	supported
ODA78969.1	c	NA	NA	NA (NA)	NA	NA	NA	NA (NA)	NA	6	6	42 (2)	supported
ODA78969.1	c	4	4	23 (23)	supported	supported	4	29 (29)	supported	4	4	42 (41)	supported
ODA78969.1	c	7	6	23 (0)	supported	supported	6	29 (0)	supported	7	6	42 (3)	supported
ODA78969.1	c	4	4	23 (22)	supported	supported	4	29 (29)	supported	4	4	42 (36)	supported

Header descriptions

#gene Gene involved, refers to Swe2 proteinGenebank ID
#reference nucleotide of the homopolymer
#Swe2 Homopolymer length in Swe2 (reference for Swe1 and Swe3)
#Swe1 Homopolymer length in Swe1
#Swe1 read Total number of reads (number of reads that support the homopolymer length in Swe1) - Canu correct and trimmed reads
#Swe3 guppy remark Support of this homopolymer in the full set of reads that was used to polish the genome - Guppy 3.0.3
#Swe3 Homopolymer length in Swe3
#Swe3 read Total number of reads (number of reads that support the homopolymer length in Swe3) - Canu correct and trimmed reads
#Swe3 guppy remark Support of this homopolymer in the full set of reads that was used to polish the genome - Guppy 3.0.3
#Dan2 Homopolymer length in Dan2 (Reference for Dan1)
#Dan1 Homopolymer length in Dan1
#Dan1 read Total number of reads (number of reads that support the homopolymer length in Dan1) - Canu correct and trimmed reads
#Dan1 guppy remark Support of this homopolymer in the full set of reads that was used to polish the genome - Guppy 3.0.3

Counts were made manually, with an estimated accuracy of +/- 10%.

Supplementary methods

GenBank query:

We queried the Assembly database on the NCBI website with:

```
("minion"[Sequencing Technology] OR "nanopore"[Sequencing Technology] OR "nanopore minion"[Sequencing Technology] OR "nanopore technologies"[Sequencing Technology] OR "nanopore technology"[Sequencing Technology] OR "nanpore"[Sequencing Technology] OR "ont"[Sequencing Technology] OR "ont minion"[Sequencing Technology] OR "oxford"[Sequencing Technology] OR "oxford nanopore"[Sequencing Technology] OR "oxford nanopore minion"[Sequencing Technology] OR "oxford nanopore technologies"[Sequencing Technology] OR "oxford nanopore technology"[Sequencing Technology]) NOT ("illumina" OR "bgi" OR "pacbio" OR "pacbio rs" OR "pacbio rs ii" OR "pacbio rsii" OR "pacbio sequel" OR "pacbio smrt" OR "pacific biosciences" OR "sequel" OR "hybrid" OR "hybrid assembly") AND (latest[filter] OR "latest genbank"[filter]) AND (all[filter] NOT "derived from surveillance project"[filter]) AND (all[filter] NOT anomalous[filter])
```

The accuracy of this query depends on the assembly metadata in GenBank.

K-mer based dot-plots:

The dot-plots based on k-mers were computed with the script *repaver.r* (available at <https://gitlab.com/gringer/bioinfscripts>) The version used here was the commit 0854f3af76cb5bc5a7e0a2b4a032518f51e210fa, with the parameters *-k 11 -style dotplot*.

Search for Dan1 / Dan2 transposable element at the breakpoints:

A first analysis was conducted with the Dfam (<https://dfam.org/home>) to search for transposable elements (TE). The parameter *organism* was set to *other*. No TE was present in the neighbourhood of the different breakpoints. We also mined the two genomes for TE with transposonPSI (<http://transposonpsi.sourceforge.net/>) with default parameters. Here also, no TE was present within the 50kb surrounding each breakpoint.

PCR:

PCR was performed to test the putative chimeric nature of a Canu contig in the Swe3 assembly, with primers UpF (AACTGTGTCTAACTAGCCCG), UpR (AGGGTCCTCATAAACTTGGC), DownF (TGTATCAGGTTCCCGAATGG), DownR (CTAGGCTGGGGAATCTTCTG), SpanF (CCATCAACTTCAGCTGCTC), SpanR (CTCCTCAATCTCCCTCTCGG), ConfF (ATCGGCGACTACCTGCAC), ConfR (CGTTCCATCGTTACCACAGC). PCR reactions were run according to the GoTaq® G2 Flexi DNA polymerase instructions (Promega), with 50 ng of template DNA, 1 mM of each forward and reverse primers, in a final volume of 25 μ L. The reaction started by initial denaturation at 95°C for 2 min, followed by 30 amplification cycles (95°C for 30 sec, 60°C for 30 sec and 72°C for 30 sec), and a final elongation for 5 min at 72°C.