

1 **A chromosome-level genome assembly for the Pacific oyster (*Crassostrea***
2 ***gigas*)**

3 Carolina Peñaloza^{1*}, Alejandro P. Gutierrez^{1,2*}, Lel Eory^{1*}, Shan Wang³, Ximing
4 Guo³, Alan L. Archibald¹, Tim P. Bean¹, Ross D. Houston^{1#}

5

6 1 The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of
7 Edinburgh, Midlothian EH25 9RG, UK.

8 2 Current address: Institute of Aquaculture, Faculty of Natural Sciences, University of
9 Stirling, Stirling FK9 4LA, UK.

10 3 Haskin Shellfish Research Laboratory, Department of Marine and Coastal
11 Sciences, Rutgers University, 6959 Miller Avenue, Port Norris, NJ 08349, USA

12

13 * These authors contributed equally to this manuscript

14 # Corresponding author Ross D. Houston ross.houston@roslin.ed.ac.uk

15

16 **Abstract**

17 The Pacific oyster (*Crassostrea gigas*) is a marine bivalve species with vital roles in
18 coastal ecosystems and aquaculture globally. While extensive genomic tools are
19 available for *C. gigas*, highly contiguous reference genomes are required to support
20 both fundamental and applied research. In the current study, high coverage long and
21 short read sequence data generated on Pacific Biosciences and Illumina platforms
22 from a single female individual specimen was used to generate an initial assembly,

23 which was then scaffolded into 10 pseudo chromosomes using both Hi-C
24 sequencing and a high density SNP linkage map. The final assembly has a scaffold
25 N50 of 58.4 Mb and a contig N50 of 1.8 Mb, representing a step advance on the
26 previously published *C. gigas* assembly. The new assembly was annotated using
27 Pacific Biosciences Iso-Seq and Illumina RNA-Seq data, identifying 30K putative
28 protein coding genes, with an average of 3.9 transcripts per gene. Annotation of
29 putative repeat elements highlighted an inverse relationship with gene density, and
30 identified putative centromeres of the metacentric chromosomes. An enrichment of
31 *Helitron* rolling circle transposable elements was observed, suggesting their
32 potential role in shaping the evolution of the *C. gigas* genome. This new
33 chromosome-level assembly will be an enabling resource for genetics and genomics
34 studies to support fundamental insight into bivalve biology, as well as for genetic
35 improvement of *C. gigas* in aquaculture breeding programmes.

36

37 **Background**

38 The Pacific oyster, *Crassostrea gigas* (Thunberg, 1793) (NCBI:txid29159), also
39 referred to as *Magallana gigas* by some authors [1, 2], is a keystone ecosystem and
40 aquaculture species [3]. Although native to the Pacific coast of Northeast Asia [4], *C.*
41 *gigas* has been introduced to all continents, except Antarctica, for farming purposes
42 [5-9]. The intensive human-mediated spread of Pacific oysters was mainly catalysed
43 by the collapse of the fishery and culture of native oyster stocks due to disease, and
44 the need to supplement depleted stocks [10, 11]. Most of these initiatives had far-
45 reaching effects on the global distribution of Pacific oysters, since several self-
46 sustaining populations became established in the wild [12, 13]. As a result, *C. gigas*

47 is now one of the most highly produced aquaculture species globally, and a
48 conspicuous invasive species in many countries [14].

49 The extent of genetic and genomic resources developed for Pacific oysters are
50 unparalleled among molluscs and other lower invertebrates [15]. Hence, they are
51 often used as model organisms to represent Lophotrochozoa, a major clade showing
52 the greatest diversity of body plans among Bilaterians [16-18]. These resources have
53 also been applied to enhance aquaculture production, with early technological
54 advances in *C. gigas* focused on developing techniques to improve production
55 through ploidy manipulation [19, 20], which later allowed the creation of the first
56 tetraploid and triploid oyster stocks [21]. Advances in DNA sequencing technologies
57 led to rapid additional resource development for this species, including extensive
58 transcriptome datasets [22-26], linkage maps using microsatellite and more recently
59 SNP markers [27, 28], and medium and high density SNP arrays [29, 30]. These
60 tools have become valuable genomic resources to enhance genetic improvement of
61 production traits, such as growth and disease resistance [31, 32]. Nevertheless, a
62 key resource for enabling genetics and genomic research in a given species is a
63 high quality reference genome. Zhang, Fang [33] published the first draft reference
64 genome assembly for *C. gigas* using a whole genome shotgun sequencing approach
65 and short read Illumina sequenced data. Interrogation of the reference genome data
66 pointed to gene expansion as a likely factor explaining the adaptation of *C. gigas* to
67 challenging marine environments, a finding that has been mirrored in a number of
68 subsequent reference genome studies for bivalve shellfish (reviewed in [34]).
69 Although a major achievement, and indeed one of the first genome assemblies for a
70 molluscan species, the publicly available reference genome is highly fragmented
71 (GenBank accession number GCA_000297895.2, 26,965 contigs, contig N50 =

72 42.3 Kb). Moreover, the previous version of this assembly (GCA_000297895.1)
73 contains many misplaced and chimeric scaffolds as revealed by alignment with
74 linkage maps [28, 30]. These issues are likely to derive from a combination of both
75 biological factors, such as the high levels of genome heterozygosity, and technical
76 factors, such as the reliance on short-read sequencing available at the time [33].
77 Therefore, highly contiguous and accurate reference genome assemblies would
78 represent valuable resources for enabling genetics and genomic research in this
79 keystone species.

80 In the current study, an improved (chromosome-level) assembly was developed for
81 *C. gigas* by harnessing high coverage Pacific Biosciences (PacBio) long-read
82 sequencing (80X), alongside accurate Illumina short read data (50x). The assembly
83 was then scaffolded to chromosome level using both Hi-C sequencing and a high-
84 density SNP linkage map, and the genome was annotated based on both Illumina
85 and PacBio transcript sequencing. This improved reference genome assembly
86 represents a step forward towards improving our understanding of fundamental
87 biological and evolutionary questions, and the genetic improvement of important
88 aquaculture production traits via genomics-enabled breeding.

89

90 **Sample collection and sequencing**

91 A single female individual collected in 2017 from Guernsey Sea Farms (Guernsey,
92 UK) was used for whole-genome resequencing with the PacBio Sequel (Pacific
93 Biosciences, Menlo Park, CA, USA) and the HiSeq X (Illumina Inc.; San Diego, CA,
94 USA) platforms. High quality dsDNA was isolated from ethanol-preserved gill tissue
95 using a cetyl trimethylammonium bromide (CTAB) method. The quality of the DNA

96 extraction was verified by the NanoDrop A260/280 and 260 /230 ratios (ND-1000)
97 and a fluorescence-based electrophoresis on a 2200 TapeStation System (Agilent
98 Technologies, USA). Using this purified DNA, three different types of libraries were
99 prepared to generate the sequencing data used for the assembly of the *C. gigas*
100 genome. The first set of libraries were generated to obtain long PacBio reads and
101 develop an initial *de novo* assembly. Two SMRTbell® libraries (chemistry v3.0) were
102 prepared and sequenced by Edinburgh Genomics (University of Edinburgh, UK)
103 across 13 SMRT cells of a PacBio Sequel system. A total of ~55 Gb of raw bases
104 with an N50 length of 12,777 bp were produced (Supplementary Figure S1). Second,
105 a paired-end sequencing library of 300 bp insert size was prepared from the same
106 individual and then used for (i) sequence error correction and (ii) quality assessment
107 of the draft genome assembly. This library was produced by Edinburgh Genomics
108 using the TruSeq DNA Nano gel free library kit (Illumina) and then sequenced on a
109 HiSeq X platform (2 x 150 bp paired-end reads). Approximately 210 million short
110 reads were obtained after quality filtering (average BQ>15 over 5 bp) and adapter
111 removal with Trimmomatic v0.38 [35]. Thirdly, a Hi-C library was generated with the
112 purpose of scaffolding the assembly into large pseudo-chromosomes. Libraries were
113 prepared using the Dovetail™ Hi-C Library Preparation Kit, following the
114 manufacturer's protocol (Dovetail™ Hi-C Kit Manual v.1.03). This final library was
115 sequenced on an Illumina HiSeq X platform (2 x 150 bp), and resulted in 500 million
116 read pairs.

117 Total RNA was extracted from two additional individual oysters (also from Guernsey
118 Sea Farms, Guernsey, UK), a male and a female, from six distinct tissues (gill,
119 mantle, stomach, heart, adductor muscle and gonads (ovaries and testis)). Full-
120 length transcripts were isolated from the tissue samples using a combination of the

121 TRIzol (Invitrogen) and the RNeasy plus minikit (Qiagen) protocols, with the inclusion
122 of a DNase treatment step. RNA quality was assessed using the Nanodrop ND-1000
123 and the Agilent 2200 TapeStation instruments. RNA extracts were quantified using a
124 QubitTM RNA assay kit (Thermo Fisher, Waltham, MA, US), and then combined in
125 equimolar quantities into a single pool for sequencing. The final RNA-pool was used
126 to obtain full-length cDNA sequences using the TeloPrime Full-Length cDNA
127 Amplification Kit V2 (Lexogen). cDNA was then sequenced across three SMRT cells
128 of a PacBio Sequel platform at the Dresden-concept Genome Center DcGC
129 (Germany). A total of 178 Gb of data comprising 1.6 million transcripts with a mean
130 length of 1.3 kb were generated for gene annotation.

131

132 **Genome features**

133 Due to the differences in genome size estimates reported in the literature for *C.*
134 *gigas* [15, 33], the DNA content of the Pacific oyster genome was also estimated in
135 the current study. To this end, the average genome size was estimated for the
136 sequenced female using the k-mer method [36] and flow cytometry [37]. For the k-
137 mer based approach, quality-filtered Illumina reads (150 bp length) were used to
138 count the frequency of different *k*-mer sizes, ranging from 15 to 23, using Jellyfish
139 v2.1.3 [36]. All *k* values evaluated showed a clear bimodal distribution, with peaks
140 occurring at a read depth of 19X and 37X (Supplementary Figure S2). The k-mer
141 frequency plots obtained are characteristic of species with highly heterozygous
142 genomes [38]. From the k-mer based analysis (k-mer = 21), the *C. gigas* genome
143 size was estimated at 534 Mb. For the genome size estimation by flow cytometry,
144 Pacific oyster nuclei were isolated and stained with propidium iodide. Two species

145 were used as internal standards for the assay, fruit fly (*Drosophila melanogaster*)
146 and zebra fish (*Danio rerio*). According to flow cytometry measurements, the
147 genome size of the female oyster sequenced in the current study was estimated at
148 640 Mb. Due to the different genome size estimates obtained by the two methods,
149 the midpoint – i.e. ~590 Mb - was used to calculate the predicted sequencing yield
150 and anticipated length for *de novo* genome assembly. The heterozygosity of the
151 Pacific oyster genome was assessed with GenomeScope v2.0 [39], based on the
152 quality filtered Illumina reads. A heterozygosity rate of 3% was estimated from the
153 21-mer based assessment of the oyster genome (Supplementary Figure S3). This
154 value is higher than the 1.3% previously reported for this species [33], but is likely
155 explained by the fact that the authors used an inbred individual for genome
156 assembly, whereas in this study an outbred female was sequenced. Although high,
157 the heterozygosity value is in the range with those reported for other bivalve
158 molluscs (e.g. 2.4% in the quagga mussel [40]).

159

160 **Genome assembly**

161 The PacBio reads were first assembled into contigs using Canu v1.8 [41] at near
162 default parameters (corrected error rate = 0.045 and raw error rate = 0.300). Contigs
163 were polished with one round of Arrow [42] followed by an additional round of
164 polishing with Pilon [43], after alignment of the post-quality filtered Illumina reads
165 with Minimap2 [44]. Compared with the genome size estimate of 590 Mb, the initially
166 assembled version of the genome was approximately two times larger than
167 expected, yielding 6,368 contigs, a total length of ~1.2 Gb, and an N50 length of 0.46
168 Mb. These results can be explained by the high frequency of highly divergent

169 haplotypes in the *C. gigas* genome, a feature that has also been observed in the
170 process of creating genome assemblies for other molluscan species [45, 46]. Whilst
171 the size of the assembled sequence could indicate that the high level of
172 heterozygosity had allowed the resolution of the two haplotypes present, we sought
173 to establish a high quality pseudo-haploid genome as a reference. To assess the
174 level of duplication in the initial assembly, a BUSCO analysis was performed [47]. By
175 searching against the metazoa_odb9 database using sea hare as a reference
176 species, 791 BUSCO genes (80.9%) were found to be duplicated. To remove
177 potentially redundant contigs by retaining only one variant of a pair of divergent
178 haplotypes, two independent approaches were taken. First, the short read data were
179 used to identify and reassign putative haplotigs with the Purge Haplotigs pipeline (-l
180 5, -m 38, -h 90) [48]. Secondly, an all-versus-all contig mapping was performed on
181 the repeat masked assembly with minimap2 v.2.2.15 [44]. Contigs were ordered
182 based on their length and matching contigs which mapped at least 30% of their
183 length and longer than 10kbp were removed as potential haplotigs. The reference
184 sequence and the mapping sequences were all removed before the next iteration.
185 The lists of curated contigs obtained independently from both methods were
186 compared and the common contigs then selected for an additional round of haplotig
187 purging. This approach resulted in a significant reduction in the number of contigs to
188 1,235, which were retained for scaffolding.

189

190 **Chromosome-level assembly using Hi-C and linkage map data**

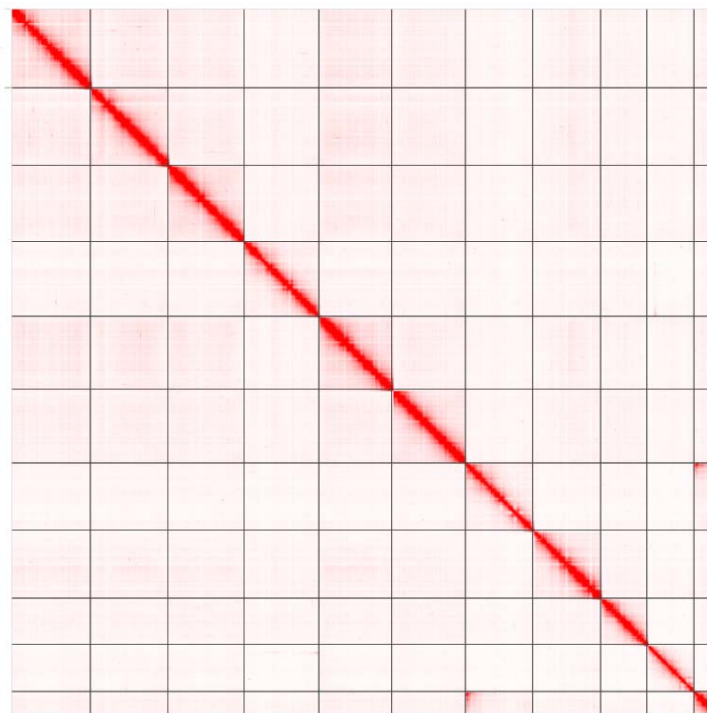
191 To generate a chromosome level assembly for *C. gigas*, Hi-C proximity ligation [49]
192 data were used to order and orientate the contigs along chromosomes. The

193 scaffolding process was carried out by Dovetail Genomics (Santa Cruz, CA, USA)
194 using the Dovetail™ Hi-C library reads to connect and order the input set of contigs.
195 After scaffolding with HiRise v2.1.7 [50], the assembled genome sequence initially
196 comprised a total of ~ 633 Mb, with a scaffold and contig N50 of 57 and 0.7 Mb and,
197 respectively. A high fraction of the assembled sequences (>92%) was contained in
198 only 11 super-scaffolds (Figure 1). However, Pacific oysters have 10 pairs of
199 chromosomes [51]. A high-density linkage map [27] was used to anchor the super-
200 scaffolds into chromosomes. SNP probes were mapped to the reference genome
201 assembly using BWA v0.78 [52]. Of the 20,353 markers on the genetic map, 17,747
202 mapped to a chromosome-level scaffold with a MAPQ above 16. The integration of
203 genetic-linkage information enabled the anchoring of two super-scaffolds onto a
204 single linkage group (LG2), resulting in an assembly with 10 major scaffolds that
205 represent all oyster chromosomes (Figure 2). Gaps were closed with PBJelly [53]
206 and again error corrected using the short read Illumina data using Pilon [43]. From
207 the remaining set of unplaced scaffolds, regions of low sequence accuracy were
208 identified based on short read coverage, following [54]. Briefly, the median read-
209 depth per 1,000 bp (non-overlapping) windows was calculated, after GC-content
210 normalization. Scaffolds with >70% of windows showing a median coverage 2SD
211 above or below the mean were removed from the analysis. All unplaced contigs and
212 scaffolds showing significant sequence identity with the Iso-Seq data were added to
213 the primary set.

214 The final Pacific oyster assembly (GenBank accession number GCA_902806645.1)
215 contains the ten expected chromosomes and 226 unplaced scaffolds, with a total
216 N50 of 58.4 Mb and 1.8 Mb for scaffold and contig lengths, respectively (Table 1).
217 This final assembly is 647 Mb in size, with the chromosome-level scaffolds

218 represented in 589 Mb of sequence. This represents a step improvement over the
219 previous version of the *C. gigas* reference genome [33], and other oyster assemblies
220 [46]. However, it should be noted that a separate chromosome-level reference
221 genome assembly is available in GenBank (accession number GCA_011032805.1)
222 from the Institute of Oceanology, Chinese Academy of Sciences. This assembly is
223 slightly shorter at 587 Mb, has a similar scaffold N50 of 61.0 Mb, and a higher contig
224 N50 of 3.1 Mb. Future comparisons between these two high quality assemblies will
225 be important to evaluate their consistency and ensure uniform use of nomenclature
226 to describe chromosomes. Furthermore, it is expected that additional high quality
227 reference genome assemblies will become available for this species, and the
228 availability of multiple assemblies is advantageous for *C. gigas* as a species with
229 high levels of intra- and inter-population genetic diversity [15]. To aid with the
230 coordination of this assembly with existing and future assemblies, the ten large
231 scaffolds in the current assembly were aligned with the karyotype using FISH probes
232 corresponding to BAC clones (Supplementary Note A). The correspondence
233 between the nomenclature of the linkage groups and scaffolds assembled in the
234 current study, and the chromosome number of the karyotype is given in
235 Supplementary Table 1. This information should enable consistency in nomenclature
236 when describing multiple genome assemblies for this species in the future.

237

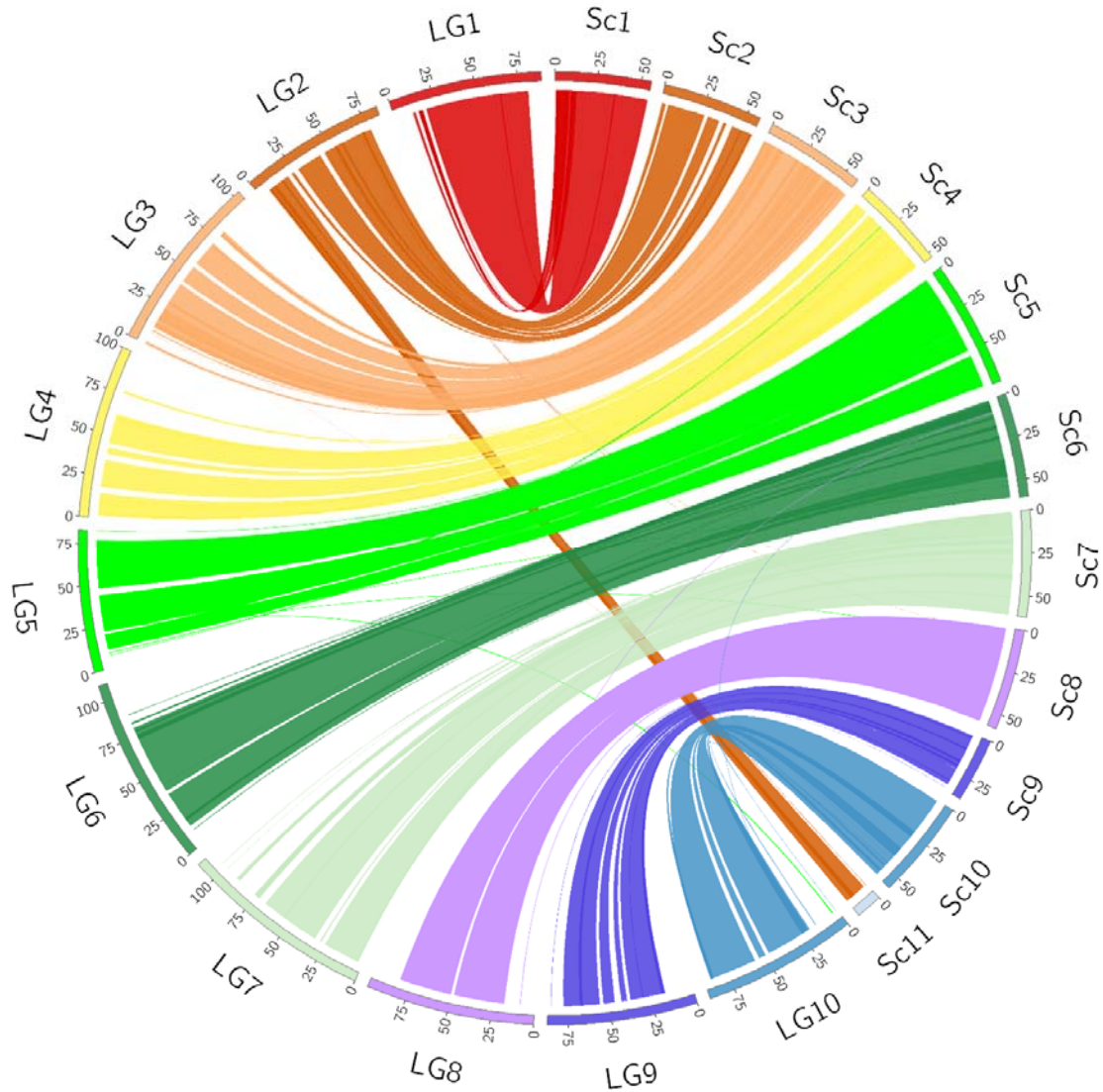


238

239

240 **Figure 1. Hi-C interaction analysis depicting the 11 super-scaffolds obtained**
241 **after using the HiRise™ scaffolding software.** Contact map is visualized using
242 Juicebox v1.11.08 [55].

243



244

245

246 **Figure 2. A high concordance between the chromosome-level scaffolds and a**
247 **high-density linkage map allowed the anchoring of two scaffolds (Sc2 and**
248 **Sc11) to a single linkage group 2 (LG2).** Ticks in each linkage group or scaffold
249 indicate lengths in 25 Mb. Scaffold (Sc) unit lengths are in Mb. Linkage group (LG)
250 units of distance are expressed in cM. Plot generated using Circos v0.69-8 [56].

251

252 **Quality assessment of reference genome**

253 Firstly, the assembled *C. gigas* genome was screened for contaminant DNA with
254 BlobTools v1 [57]. All scaffolds and contigs had a top hit match to *C. gigas*
255 (Supplementary Figure S5). Second, to assess the completeness of the assembled
256 genome, a BUSCO analysis was performed. From the curated list of single copy
257 genes, 934 (95.5%) were found in the assembly, of which 917 (93.8%) were single-
258 copy genes, 17 (1.7%) were duplicated and 38 (3.9%) were missing. Finally, to
259 evaluate the accuracy of the reconstructed *C. gigas* genome, structural variants were
260 called with Sniffles [58], after alignment of the PacBio raw reads with ngmlr v0.2.7.
261 Variants with a minimum size of 50 bp for which the ratio of high quality reads for the
262 assembly (reference) variant was below 0.2 were considered assembly errors (Table
263 S2).

264

265

266

267

268

269

270

271

272

273

274 **Table 1.** Genome assembly statistics and annotation of *C. gigas*

Genome assembly	
A) Genome	
GC content	33.25%
Total size (bp)	647,887,097
Contigs	
N° contigs	711
N50 length (bp)	1,813,842
Longest (bp)	11,935,632
Scaffolds	
N° scaffolds	236
N50 length (bp)	58,462,999
Longest (bp)	73,550,375
B) Genome annotation	
N° Transposable elements	
LTR	22,828
LINE	41,781
DNA	634,611
Total	699,220
Protein coding genes	

N°	30,844
Mean transcript length (bp)	7,527
Mean coding sequence length (bp)	1,175
Mean exon length (bp)	248
<i>Functional annotation</i>	
GO	18,750
KO	11,390

275

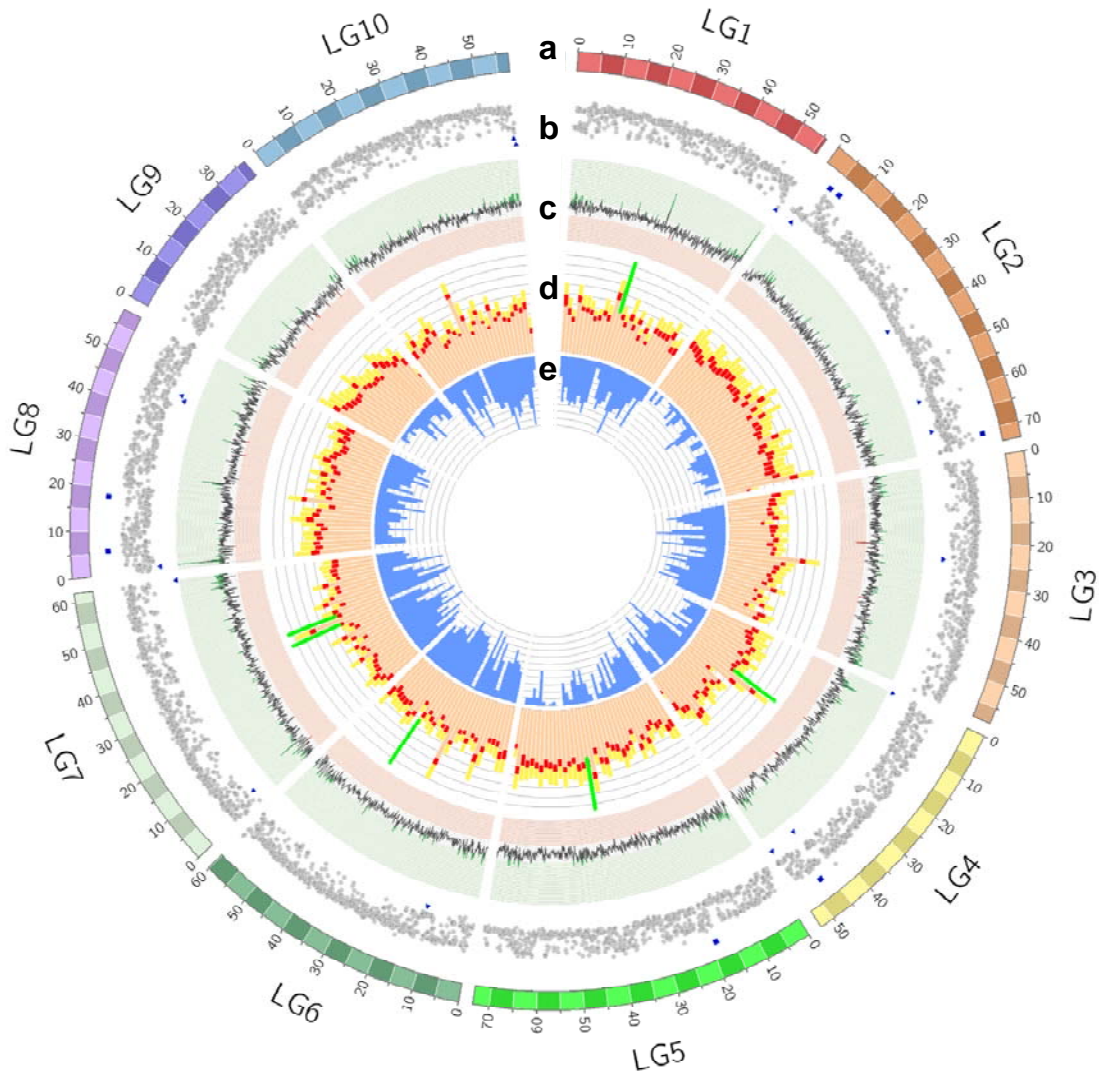
276

277 **Genome annotation**

278 Genome annotation was carried out using long-read PacBio Iso-Seq data from six
279 tissues and the Illumina short-read RNA-Seq data from Zhang [33]. Short-read data
280 was mapped to the reference assembly with STAR v.2.5.1b [59]. Transcript models
281 were created by BRAKER v.2.1.5 [60] using only the paired-end RNA-seq datasets
282 (see Supplementary Table S3). Multi exon transcripts that were expressed in at least
283 two tissues at an expression level over 1TPM were retained. Iso-Seq raw sub-reads
284 were processed with SMRT Link v7.0 (Pacific Biosciences) to obtain Circular
285 Consensus Sequences (CCS) using a '--min-rq of 0.9'. The Iso-Seq CCS reads were
286 mapped with minimap2 v.2.16 [44] and the transcript models were called using the
287 TAMA package [61] (see Supplementary Note B). Protein-coding transcripts and
288 translation start and end positions were predicted by mapping known protein
289 sequences from UniRef90 [62] to the oyster transcripts by Diamond v.0.9.31 [63].
290 Those models that contained a frameshift within the coding sequence were classified
291 as pseudogenes.

292 The final annotation of the assembled *C. gigas* genome contains 35,527 genes, of
293 which 30,844 are protein coding, 4,001 represent non-coding RNA genes and 682
294 were classified as pseudogenes. Among the protein coding genes, 14,293 (49%)
295 contained putative alternative spliced transcripts, with an average of 3.9 transcripts
296 per gene. The gene models predicted for *C. gigas* were functionally annotated using
297 the Blast2GO pipeline [64], and KEGG orthology (KO) groups were assigned using
298 KOBAS v2.0 [65]. Approximately, 18,750 (61%) of the predicted protein coding
299 genes were assigned functional labels (Table 1). This reference genome assembly
300 has also been annotated by the NCBI annotation team, who used the extensive short
301 read transcriptome data available for *C. gigas* to annotate 38,296 genes (31,371
302 protein coding, 6,837 non-coding, 88 pseudogenes) and a total of 73,946 transcripts
303 [66].

304



305

306 **Figure 3. Circos plot depicting genome features across the 10 oyster**

307 **chromosomes. (a) Oyster chromosomes (LG1–LG10 on a Mb scale). (b) short-read**

308 **coverage plot. Coverage within 2SD of the mean are shown as grey circles.**

309 **Abnormal sequence coverage ($\pm 2SD$ from the mean) are indicated with a blue**

310 **square or triangle, respectively. (c) GC content percentage (>35% in green; <31% in**

311 **red). (d) Distribution of repeat elements: DNA transposons (light orange bar),**

312 **retrotransposon TEs (red bar) and novel repeat elements (yellow bar). The location**

313 **of centromeres is indicated with a green line. (e) Gene density (range 50-150). For**

314 tracks (b) and (c) a window size of 0.1 Mb was used, whereas for tracks (d) and (e)
315 the size was increased to 0.2 Mb.

316

317 **Repeat element annotation**

318 Known Pacific oyster specific repeat sequences were identified in the genome
319 assembly using RepeatMasker v.4.0.7 [67] with a combined repeat database
320 (Dfam_Consensus-20170127 and RepBase-20170127) [68, 69] with parameters '-s -
321 species "Crassostrea gigas" -e ncbi'. Besides the 972 repeat families contained in
322 the RepeatMasker library an additional 1,827 novel repeat families were identified by
323 RepeatModeler v.1.0.11 [70]. This novel repeat library was used to identify the
324 location of novel elements in the newly built assembly. For comparison, an exact
325 same search was performed on the older version of the *C. gigas* genome assembly
326 (GenBank assembly accession GCA_000297895.2).

327 Overall, a higher number of repetitive elements were identified in our assembly
328 compared to the previous genome assembly (Figure S6). Nearly 43% of the Pacific
329 oyster genome was represented by repeat elements. Repetitive sequences were
330 distributed unevenly along the *C. gigas* chromosomes. In general, an inverse
331 relationship between the total number of repeat elements and gene density was
332 observed (Figure 3 d-e). Among the different classes of repeat elements, significant
333 negative correlations were found between gene density and (i) retrotransposons of
334 the LTR type (corr = -0.61; $P = 2.2 \times 10^{-16}$), (ii) Non-LTR retrotransposons (corr = -
335 0.28; $P = 5.4 \times 10^{-7}$), (iii) satellite DNA (corr = -0.29; $P = 4.5 \times 10^{-7}$), (iv) simple
336 repeats (corr = -0.33; $P = 4.7 \times 10^{-9}$), and (v) DNA transposons (corr = -0.59; $P = 2.2$
337 $\times 10^{-16}$). The centromere of five metacentric chromosomes were located after

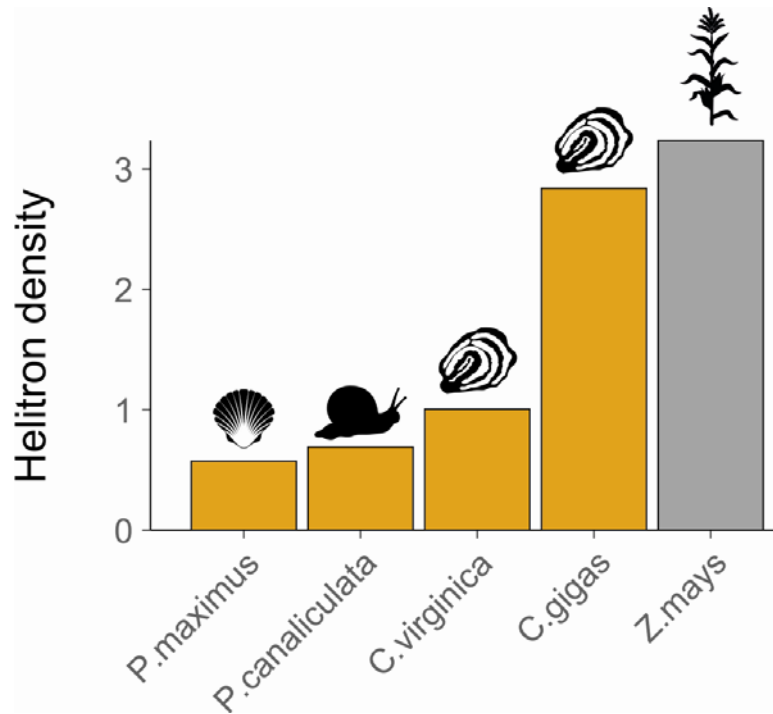
338 aligning six centromere-associated microsatellite markers to the assembly [71]
339 (Table S4). Four of these five centromere regions co-localize with genomic windows
340 enriched for repetitive elements (Figure 3d). Among repetitive elements,
341 transposable elements (TEs) were the most common, and accounted for 36% of the
342 assembled genome. Consistent with previous studies [46], the oyster genome is
343 dominated by DNA transposons (32 % of the genome assembly) (Table 1), with
344 *Helitrons* being the most abundant superfamily (Supplementary Figures S7-8).

345

346 **Characterization of *Helitrons* in the Pacific oyster genome**

347 *Helitrons* are rolling-circle transposable elements that have the ability to capture host
348 gene fragments [72]. In maize, *Helitrons* have significantly influenced genome
349 evolution, leading to genome variation among lines [73] and to a notable
350 diversification of transcripts via exon shuffling of thousands of genes [74]. To refine
351 the annotation of Pacific oyster *Helitrons*, a structure-based search [75] was
352 performed in addition to the homology based approach described above. The
353 localization of these elements was heterogeneous across the Pacific oyster
354 chromosomes, with LG5 and LG8 showing a higher density of elements (>1SD
355 above the average across chromosomes) (Figure S9). *Helitrons* in plant and animal
356 genomes tend to accumulate in gene-poor regions [76]. However, this bias is less
357 evident in *C. gigas*, with no significant association found between gene density and
358 the number of *Helitrons* within a region. A comparison with other molluscan
359 reference genome assemblies revealed that *C. gigas* had a remarkably high number
360 of predicted *Helitron* related sequences (Figure 4).

361



362

363 **Figure 4. Density of *Helitrons* identified across four molluscan genomes**
364 **(orange bars), including maize as a reference species (grey bar).** The reference
365 genome assembled for *C. gigas* was compared to the king scallop (*Pecten maximus*;
366 GCF_902652985.1), golden apple snail (*Pomacea canaliculata*; GCF_003073045.1),
367 and Atlantic oyster (*Crassostrea virginica*; GCF_002022765.2), with maize included
368 as a reference species (*Zea mays*; GCF_000005005.2). *Helitron* density is
369 expressed as the number of conserved 3' ends over genome size (in Mb).

370

371 The Pacific oyster *Helitron*-like sequences possess the basic expected structure as
372 observed in other taxa: TC sequence at the 5' termini, CTAG motif on the 3'-
373 terminus, and a 16-20 bp palindromic sequence that can form a hairpin structure
374 upstream of the 3'-end. Likewise, they were also found to preferentially insert (86%
375 of the cases) between the 5'-A and 3'-T nucleotides of the host AT target sites. Of
376 the 751 intact *Helitrons* discovered through the *in silico* screening, 627 elements had

377 a high 3'-end pairwise sequence similarity (identity of ~85 % over 30 bp), suggesting
378 they belong to the same family [76]. Notably, a significant fraction of these elements
379 (261 out of 751) had sub-terminal inverted repeats, as revealed by a screening of
380 their paired terminal ends with the Inverted Repeats Database (IRDB;
381 <https://tandem.bu.edu/cgi-bin/irdb/irdb.exe>). This structural feature is characteristic of
382 an alternative variant of *Helitrons* called *Helentrons*, which in its non-autonomous
383 form known as HINE (Helentron-associated INterspersed Elements) has been
384 recently associated to large numbers of satellite DNA in the oyster genome [77].

385 *Helitrons* have been observed to capture gene fragments in species such as maize
386 and the little brown bat (*Myotis lucifugus*) [78, 79]. In *C. gigas*, a BLASTX [80] search
387 against the UniRef database revealed that only 17 *Helitrons* (2%) carried gene
388 fragments; alignment lengths >50 with at least 85% identity were considered a
389 match. The Pacific oyster *Helitron*-like sequences were relatively short (mean = 1092
390 bp; SD = 557 bp), and lacked the main enzymatic hallmarks of autonomous
391 elements (i.e., REP/Helicase domains). Non-autonomous *Helitrons* require the
392 transposase expressed by their autonomous counterparts in order to amplify. Due to
393 the fact this study did not detect evidence for the presence of autonomous mobile
394 sequences in the Pacific oyster genome, these abundant *Helitron* elements are likely
395 to be inactive, suggesting they are remnants of high levels of past activity in the
396 evolutionary history of *C. gigas*.

397

398

399

400

401 **Conclusion**

402 The new chromosome-level *C. gigas* genome assembly presented herein has a
403 scaffold N50 of 58.4 Mb and a contig N50 of 1.8 Mb, representing a step advance on
404 the previously published assembly, and will complement other high quality
405 assemblies available or becoming available in the near-future. Approximately 30K
406 putative protein coding genes were identified with an average of 3.9 transcripts per
407 gene. DNA transposons dominated the repeat elements detected in the assembly,
408 with *Helitrons* being found at a level substantially higher level than other molluscan
409 species, suggesting their potential role in shaping the evolution of the *C. gigas*
410 genome. The availability of a chromosome-level genome assembly is expected to
411 support applied and fundamental research in this keystone ecological and
412 aquaculture species.

413

414 **Availability of supporting data**

415 Raw sequencing data has been submitted to the European Nucleotide Archive
416 (ENA) under study accession number PRJEB35351. The genomic short read data
417 are under accessions numbers ERX3728455, ERX3728453, ERX3728482,
418 ERX3728546, ERX3728630 and ERX3728636; the raw reads of the Hi-C library are
419 under accession numbers ERX3722775. PacBio Iso-Seq reads of pooled samples
420 are available under accession numbers ERX3721883, ERX3722678 and
421 ERX3722679. Raw PacBio reads from the nuclear DNA are available under
422 accessions ERX3761471, ERX3761586, ERX3761587, ERX3761621, ERX3761714,
423 ERX3761715, ERX3761720, ERX3762151, ERX3762342, ERX3762370,

424 ERX3762371, ERX3762372 and ERX3762598. The Pacific oyster genome assembly
425 is available at GenBank under accession number GCA_902806645.1.

426

427 **Abbreviations**

428 bp: base pairs; BQ: base quality; BUSCO: Benchmarking Universal Single-Copy
429 Orthologs; cM: centimorgan; cDNA: coding DNA ;DNA: deoxyribonucleic acid; Gb:
430 giga base pairs ; GC: guanine-cytosine; Gb: gigabase pairs; kb: kilobase pairs;
431 KEGG: Kyoto encyclopedia of genes and genomes; MAPQ: mapping quality; Mb:
432 megabase pairs; N50: median size; PacBio: Pacific Biosciences; RNA: ribonucleic
433 acid; RNA-Seq: RNA-sequencing; SMRT: single-molecule real-time.

434

435 **Acknowledgements**

436 The authors thank Katy Monteith, Darren Obbard and Carl Tucker for providing
437 controls for the flow cytometry assay; Guernsey Sea Farms for rearing the female
438 oyster used for generating the assembly; and Manu Kumar Gundappa and Richard
439 Kuo from the Roslin Institute for their technical advice during the assembly and
440 annotation steps. This work was supported by funding from the Natural Environment
441 Research Council (NE/P010695/1) and Biotechnology and Biological Sciences
442 Research Council (BB/S004343/1, BB/P013759/1 and BB/P013740/1). The
443 cytogenetic mapping of BACs was supported by a grant from U.S. Department of
444 Agriculture (2009-35205-05052).

445

446 **References**

- 447 1. Salvi D, Macali A and Mariottini P. Molecular phylogenetics and systematics of the bivalve
448 family Ostreidae based on rRNA sequence-structure models and multilocus species tree.
449 PloS one. 2014;9 9:e108696-e. doi:10.1371/journal.pone.0108696.
- 450 2. Salvi D and Mariottini P. Molecular taxonomy in 2D: a novel ITS2 rRNA sequence-structure
451 approach guides the description of the oysters' subfamily Saccostreinae and the genus
452 Magallana (Bivalvia: Ostreidae). Zoological Journal of the Linnean Society. 2016;179 2:263-
453 76. doi:10.1111/zoj.12455.
- 454 3. FAO. 2020. Rome, Italy.
- 455 4. Wang H, Qian L, Liu X, Zhang G and Guo X. Classification of a Common Cupped Oyster from
456 Southern China. Journal of Shellfish Research. 2010;29:857-66. doi:10.2983/035.029.0420.
- 457 5. Robinson T, Griffiths C, Tonin A, Bloomer P and Hare M. Naturalized populations of oysters,
458 *Crassostrea gigas* along the South African coast: Distribution, abundance and population
459 structure. Journal of Shellfish Research. 2009;24:443-50. doi:10.2983/0730-
460 8000(2005)24[443:NPOOCG]2.0.CO;2.
- 461 6. Anglès d'Auriac MB, Rinde E, Norling P, Lapègue S, Staalstrøm A, Hjermand DØ, et al. Rapid
462 expansion of the invasive oyster *Crassostrea gigas* at its northern distribution limit in
463 Europe: Naturally dispersed or introduced? PLOS ONE. 2017;12 5:e0177481.
464 doi:10.1371/journal.pone.0177481.
- 465 7. Carrasco MF and Barón PJ. Analysis of the potential geographic range of the Pacific oyster
466 *Crassostrea gigas* (Thunberg, 1793) based on surface seawater temperature satellite data
467 and climate charts: the coast of South America as a study case. Biological Invasions. 2010;12
468 8:2597-607. doi:10.1007/s10530-009-9668-0.
- 469 8. Miller PA, Elliott NG, Koutoulis A, Kube PD and Vaillancourt RE. Genetic Diversity of Cultured,
470 Naturalized, and Native Pacific Oysters, *Crassostrea Gigas*, Determined from Multiplexed
471 Microsatellite Markers. Journal of Shellfish Research. 2012;31 3:611-7, 7.
- 472 9. Meistertzheim A-L, Arnaud-Haond S, Boudry P and Thébault M-T. Genetic structure of wild
473 European populations of the invasive Pacific oyster *Crassostrea gigas* due to aquaculture
474 practices. Marine Biology. 2013;160 2:453-63. doi:10.1007/s00227-012-2102-7.
- 475 10. Shatkin G, Shumway S and Hawes R. Considerations regarding the possible introduction of
476 the Pacific oyster, *Crassostrea gigas*, to the Gulf of Maine: a review of global experience.
477 Journal of Shellfish Research. 1997;16:463-78.
- 478 11. Jones MC, Dye SR, Pinnegar Jk, Warren R and Cheung WWL. Applying distribution model
479 projections for an uncertain future: the case of the Pacific oyster in UK waters. Aquatic
480 Conservation: Marine and Freshwater Ecosystems. 2013;23 5:710-22. doi:10.1002/aqc.2364.
- 481 12. Wrange A-L, Valero J, Harkestad LS, Strand Ø, Lindegarth S, Christensen HT, et al. Massive
482 settlements of the Pacific oyster, *Crassostrea gigas*, in Scandinavia. Biological Invasions.
483 2010;12 5:1145-52. doi:10.1007/s10530-009-9535-z.
- 484 13. Herbert RJH, Humphreys J, Davies CJ, Roberts C, Fletcher S and Crowe TP. Ecological impacts
485 of non-native Pacific oysters (*Crassostrea gigas*) and management measures for protected
486 areas in Europe. Biodiversity and Conservation. 2016;25 14:2835-65. doi:10.1007/s10531-
487 016-1209-4.
- 488 14. Miossec L, Le Deuff R-M and Gouletquer P. Alien Species Alert: *Crassostrea gigas* (Pacific
489 Oyster). ICES Cooperative Research Report. 2009;229:42.
- 490 15. Hedgecock D, Gaffney PM, Gouletquer P, Guo X, Reece K and Warr GW. The case for
491 sequencing the Pacific oyster genome. Journal of Shellfish Research. 2005;24 2:429-41, 13.
- 492 16. Schwartz J, Réalis-Doyelle E, Dubos M-P, Lefranc B, Leprince J and Favrel P. Characterization
493 of an evolutionarily conserved calcitonin signalling system in a lophotrochozoan, the Pacific
494 oyster (*Crassostrea gigas*). The Journal of Experimental Biology. 2019;222 13:jeb201319.
495 doi:10.1242/jeb.201319.

- 496 17. Lafont M, Petton B, Vergnes A, Pauletto M, Segarra A, Gourbal B, et al. Long-lasting antiviral
497 innate immune priming in the Lophotrochozoan Pacific oyster, *Crassostrea gigas*. Scientific
498 Reports. 2017;7 1:13143. doi:10.1038/s41598-017-13564-0.
- 499 18. Kocot KM. On 20 years of Lophotrochozoa. Organisms Diversity & Evolution. 2016;16 2:329-
500 43. doi:10.1007/s13127-015-0261-3.
- 501 19. Allen SK and Downing SL. Performance of triploid Pacific oysters, *Crassostrea gigas*
502 (Thunberg). I. Survival, growth, glycogen content, and sexual maturation in yearlings. Journal
503 of Experimental Marine Biology and Ecology. 1986;102 2:197-208.
504 doi:[https://doi.org/10.1016/0022-0981\(86\)90176-0](https://doi.org/10.1016/0022-0981(86)90176-0).
- 505 20. Downing SL and Allen SK. Induced triploidy in the Pacific oyster, *Crassostrea gigas*: Optimal
506 treatments with cytochalasin B depend on temperature. Aquaculture. 1987;61 1:1-15.
507 doi:[https://doi.org/10.1016/0044-8486\(87\)90332-2](https://doi.org/10.1016/0044-8486(87)90332-2).
- 508 21. Guo X, DeBrosse GA and Allen SK. All-triploid Pacific oysters (*Crassostrea gigas* Thunberg)
509 produced by mating tetraploids and diploids. Aquaculture. 1996;142 3:149-61.
510 doi:[https://doi.org/10.1016/0044-8486\(95\)01243-5](https://doi.org/10.1016/0044-8486(95)01243-5).
- 511 22. Riviere G, Klopp C, Ibouniyamine N, Huvet A, Boudry P and Favrel P. GigaTON: an extensive
512 publicly searchable database providing a new reference transcriptome in the pacific oyster
513 *Crassostrea gigas*. BMC Bioinformatics. 2015;16 1:401. doi:10.1186/s12859-015-0833-4.
- 514 23. Kim B-M, Kim K, Choi I-Y and Rhee J-S. Transcriptome response of the Pacific oyster,
515 *Crassostrea gigas* susceptible to thermal stress: A comparison with the response of tolerant
516 oyster. Molecular & Cellular Toxicology. 2017;13 1:105-13. doi:10.1007/s13273-017-0011-z.
- 517 24. Yue C, Li Q and Yu H. Gonad Transcriptome Analysis of the Pacific Oyster *Crassostrea gigas*
518 Identifies Potential Genes Regulating the Sex Determination and Differentiation Process.
519 Marine biotechnology (New York, NY). 2018;20 2:206-19. doi:10.1007/s10126-018-9798-4.
- 520 25. Feng D, Li Q, Yu H, Zhao X and Kong L. Comparative Transcriptome Analysis of the Pacific
521 Oyster *Crassostrea gigas* Characterized by Shell Colors: Identification of Genetic Bases
522 Potentially Involved in Pigmentation. PLOS ONE. 2015;10 12:e0145257.
523 doi:10.1371/journal.pone.0145257.
- 524 26. Zhang F, Hu B, Fu H, Jiao Z, Li Q and Liu S. Comparative Transcriptome Analysis Reveals
525 Molecular Basis Underlying Fast Growth of the Selectively Bred Pacific Oyster, *Crassostrea*
526 *gigas*. Frontiers in Genetics. 2019;10 610 doi:10.3389/fgene.2019.00610.
- 527 27. Gutierrez AP, Bean TP, Hooper C, Stenton CA, Sanders MB, Paley RK, et al. A Genome-Wide
528 Association Study for Host Resistance to Ostreid Herpesvirus in Pacific Oysters (*Crassostrea*
529 *gigas*). G3: Genes|Genomes|Genetics. 2018;8 4:1273-80. doi:10.1534/g3.118.200113.
- 530 28. Hedgecock D, Shin G, Gracey AY, Den Berg DV and Samanta MP. Second-Generation Linkage
531 Maps for the Pacific Oyster *Crassostrea gigas* Reveal Errors in Assembly of Genome
532 Scaffolds. G3 (Bethesda, Md). 2015;5 10:2007-19. doi:10.1534/g3.115.019570.
- 533 29. Qi H, Song K, Li C, Wang W, Li B, Li L, et al. Construction and evaluation of a high-density SNP
534 array for the Pacific oyster (*Crassostrea gigas*). PLOS ONE. 2017;12 3:e0174007.
535 doi:10.1371/journal.pone.0174007.
- 536 30. Gutierrez AP, Turner F, Gharbi K, Talbot R, Lowe NR, Peñaloza C, et al. Development of a
537 Medium Density Combined-Species SNP Array for Pacific and European Oysters (*Crassostrea*
538 *gigas* and *Ostrea edulis*). G3 (Bethesda). 2017;7 7:2209-18. doi:10.1534/g3.117.041780.
- 539 31. Gutierrez AP, Matika O, Bean TP and Houston RD. Genomic Selection for Growth Traits in
540 Pacific Oyster (*Crassostrea gigas*): Potential of Low-Density Marker Panels for Breeding
541 Value Prediction. Frontiers in Genetics. 2018;9 391 doi:10.3389/fgene.2018.00391.
- 542 32. Gutierrez AP, Symonds J, King N, Steiner K, Bean TP and Houston RD. Potential of genomic
543 selection for improvement of resistance to ostreid herpesvirus in Pacific oyster (*Crassostrea*
544 *gigas*). Animal Genetics. 2020;51 2:249-57. doi:10.1111/age.12909.
- 545 33. Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome reveals stress adaptation
546 and complexity of shell formation. Nature. 2012;490 7418:49-54. doi:10.1038/nature11413.

- 547 34. Gomes-dos-Santos A, Lopes-Lima M, Castro LFC and Froufe E. Molluscan genomics: the road
548 so far and the way forward. *Hydrobiologia*. 2020;847 7:1705-26. doi:10.1007/s10750-019-
549 04111-1.
- 550 35. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
551 data. *Bioinformatics*. 2014;30 15:2114-20. doi:10.1093/bioinformatics/btu170.
- 552 36. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel counting of
553 occurrences of k-mers. *Bioinformatics*. 2011;27 6:764-70.
554 doi:10.1093/bioinformatics/btr011.
- 555 37. Dolezel J and Bartos J. Plant DNA Flow Cytometry and Estimation of Nuclear Genome Size.
556 *Annals of Botany*. 2005;95 1:99-110. doi:10.1093/aob/mci005.
- 557 38. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient de novo
558 assembly of highly heterozygous genomes from whole-genome shotgun short reads.
559 *Genome research*. 2014;24 doi:10.1101/gr.170720.113.
- 560 39. Ranallo-Benavidez TR, Jaron KS and Schatz MC. GenomeScope 2.0 and Smudgeplot for
561 reference-free profiling of polyploid genomes. *Nature Communications*. 2020;11 1:1432.
562 doi:10.1038/s41467-020-14998-3.
- 563 40. Calcino AD, de Oliveira AL, Simakov O, Schwaha T, Zieger E, Wollesen T, et al. The quagga
564 mussel genome and the evolution of freshwater tolerance. *DNA Research*. 2019;26 5:411-22.
565 doi:10.1093/dnares/dsz019.
- 566 41. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH and Phillippy AM. Canu: scalable and
567 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome
568 Research*. 2017; doi:10.1101/gr.215087.116.
- 569 42. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid, finished
570 microbial genome assemblies from long-read SMRT sequencing data. *Nature Methods*.
571 2013;10 6:563-9. doi:10.1038/nmeth.2474.
- 572 43. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated
573 tool for comprehensive microbial variant detection and genome assembly improvement.
574 *PloS one*. 2014;9 11:e112963-e. doi:10.1371/journal.pone.0112963.
- 575 44. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34
576 18:3094-100. doi:10.1093/bioinformatics/bty191.
- 577 45. Takeuchi T, Kawashima T, Koyanagi R, Gyoja F, Tanaka M, Ikuta T, et al. Draft Genome of the
578 Pearl Oyster *Pinctada fucata*: A Platform for Understanding Bivalve Biology. *DNA Research*.
579 2012;19 2:117-30. doi:10.1093/dnares/dss005.
- 580 46. Wang X, Xu W, Wei L, Zhu C, He C, Song H, et al. Nanopore Sequencing and De Novo
581 Assembly of a Black-Shell Pacific Oyster (*Crassostrea gigas*) Genome. *Frontiers in
582 Genetics*. 2019;10 1211 doi:10.3389/fgene.2019.01211.
- 583 47. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing
584 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.
585 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.
- 586 48. Roach MJ, Schmidt SA and Borneman AR. Purge Haplotigs: allelic contig reassignment for
587 third-gen diploid genome assemblies. *BMC Bioinformatics*. 2018;19 1:460.
588 doi:10.1186/s12859-018-2485-7.
- 589 49. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, et al.
590 Comprehensive mapping of long-range interactions reveals folding principles of the human
591 genome. *Science (New York, NY)*. 2009;326 5950:289-93. doi:10.1126/science.1181369.
- 592 50. Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, et al. Chromosome-scale
593 shotgun assembly using an in vitro method for long-range linkage. *Genome research*.
594 2016;26 3:342-50. doi:10.1101/gr.193474.115.
- 595 51. Thiriot-Quievreux C. Review of the literature on bivalve cytogenetics in the last ten years.
596 *Cahiers de Biologie Marine*. 2002;43:17-26.

- 597 52. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
598 Bioinformatics (Oxford, England). 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.
- 599 53. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the Gap: Upgrading Genomes
600 with Pacific Biosciences RS Long-Read Sequencing Technology. PLOS ONE. 2012;7 11:e47768.
601 doi:10.1371/journal.pone.0047768.
- 602 54. Warr A, Robert C, Hume D, Archibald AL, Deeb N and Watson M. Identification of Low-
603 Confidence Regions in the Pig Reference Genome (Sscrofa10.2). Frontiers in Genetics.
604 2015;6 338 doi:10.3389/fgene.2015.00338.
- 605 55. Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, et al. Juicebox
606 Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom. Cell Syst.
607 2016;3 1:99-101. doi:10.1016/j.cels.2015.07.012.
- 608 56. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: an
609 information aesthetic for comparative genomics. Genome research. 2009;19 9:1639-45.
610 doi:10.1101/gr.092759.109.
- 611 57. Laetsch D and Blaxter M. BlobTools: Interrogation of genome assemblies. F1000Res, 2017.
- 612 58. Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, et al. Accurate
613 detection of complex structural variations using single-molecule sequencing. Nature
614 Methods. 2018;15 6:461-8. doi:10.1038/s41592-018-0001-7.
- 615 59. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal
616 RNA-seq aligner. Bioinformatics. 2012;29 1:15-21. doi:10.1093/bioinformatics/bts635.
- 617 60. Hoff KJ, Lomsadze A, Borodovsky M and Stanke M. Whole-Genome Annotation with BRAKER.
618 Methods in molecular biology (Clifton, NJ). 2019;1962:65-95. doi:10.1007/978-1-4939-9173-
619 0_5.
- 620 61. Kuo RI, Cheng Y, Smith J, Archibald AL and Burt DW. Illuminating the dark side of the human
621 transcriptome with TAMA Iso-Seq analysis. bioRxiv. 2019:780015. doi:10.1101/780015.
- 622 62. Suzek BE, Wang Y, Huang H, McGarvey PB and Wu CH. UniRef clusters: a comprehensive and
623 scalable alternative for improving sequence similarity searches. Bioinformatics. 2015;31
624 6:926-32. doi:10.1093/bioinformatics/btu739.
- 625 63. Buchfink B, Xie C and Huson DH. Fast and sensitive protein alignment using DIAMOND.
626 Nature Methods. 2015;12 1:59-60. doi:10.1038/nmeth.3176.
- 627 64. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M and Robles M. Blast2GO: a universal
628 tool for annotation, visualization and analysis in functional genomics research.
629 Bioinformatics. 2005;21 18:3674-6. doi:10.1093/bioinformatics/bti610.
- 630 65. Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, et al. KOBAS 2.0: a web server for annotation
631 and identification of enriched pathways and diseases. Nucleic Acids Research. 2011;39
632 suppl_2:W316-W22. doi:10.1093/nar/gkr483.
- 633 66. https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Crassostrea_gigas/102/ Accessed
634 01 September 2020.
- 635 67. RepeatMasker: <http://www.repeatmasker.org>. Accessed 17 April 2020.
- 636 68. Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, et al. The Dfam database of
637 repetitive DNA families. Nucleic Acids Res. 2016;44 D1:D81-D9. doi:10.1093/nar/gkv1272.
- 638 69. Bao W, Kojima KK and Kohany O. Repbase Update, a database of repetitive elements in
639 eukaryotic genomes. Mobile DNA. 2015;6 1:11. doi:10.1186/s13100-015-0041-9.
- 640 70. RepeatModeler: <http://www.repeatmasker.org>. Accessed April 17 2020.
- 641 71. Hubert S, Cognard E and Hedgecock D. Centromere mapping in triploid families of the Pacific
642 oyster *Crassostrea gigas* (Thunberg). Aquaculture. 2009;288 3:172-83.
643 doi:https://doi.org/10.1016/j.aquaculture.2008.12.006.
- 644 72. Kapitonov VV and Jurka J. Rolling-circle transposons in eukaryotes. Proceedings of the
645 National Academy of Sciences. 2001;98 15:8714-9. doi:10.1073/pnas.151269298.

- 646 73. Morgante M, Brunner S, Pea G, Fengler K, Zuccolo A and Rafalski A. Gene duplication and
647 exon shuffling by helitron-like transposons generate intraspecies diversity in maize. *Nature*
648 *Genetics*. 2005;37 9:997-1002. doi:10.1038/ng1615.
- 649 74. Barbaglia AM, Klusman KM, Higgins J, Shaw JR, Hannah LC and Lal SK. Gene capture by
650 Helitron transposons reshuffles the transcriptome of maize. *Genetics*. 2012;190 3:965-75.
651 doi:10.1534/genetics.111.136176.
- 652 75. Hu K, Xu K, Wen J, Yi B, Shen J, Ma C, et al. Helitron distribution in Brassicaceae and whole
653 Genome Helitron density as a character for distinguishing plant species. *BMC bioinformatics*.
654 2019;20 1:354-. doi:10.1186/s12859-019-2945-8.
- 655 76. Yang L and Bennetzen J. Structure-based discovery and description of plant and animal
656 Helitrons. *Proceedings of the National Academy of Sciences of the United States of America*.
657 2009;106:12832-7. doi:10.1073/pnas.0905563106.
- 658 77. Vojvoda Zeljko T, Pavlek M, Meštrović N and Plohl M. Satellite DNA-like repeats are
659 dispersed throughout the genome of the Pacific oyster *Crassostrea gigas* carried by
660 Helitron non-autonomous mobile elements. *Scientific Reports*. 2020;10 1:15107.
661 doi:10.1038/s41598-020-71886-y.
- 662 78. Yang L and Bennetzen JL. Distribution, diversity, evolution, and survival of Helitrons in the
663 maize genome. *Proc Natl Acad Sci U S A*. 2009;106 47:19922-7.
664 doi:10.1073/pnas.0908008106.
- 665 79. Pritham EJ and Feschotte C. Massive amplification of rolling-circle transposons in the lineage
666 of the bat *Myotis lucifugus*. *Proceedings of the National Academy of Sciences*. 2007;104
667 6:1895-900. doi:10.1073/pnas.0609601104.
- 668 80. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and
669 PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*.
670 1997;25 17:3389-402. doi:10.1093/nar/25.17.3389.

671