1    Genome Report: A blue mussel chromosome-scale genome assembly for aquaculture, marine
2    ecology and evolution

3    Tim Regan[1], Tiago S. Hori[2], Tim P. Bean[1]

4    [1] The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of Edinburgh, UK.

5    [2] Atlantic Aqua Farms Ltd., Charlottetown, PE, Canada.

6

7    **Keywords:** Aquaculture, evolution, bivalve, mussel, *Mytilus edulis*

8

**Abstract**

The blue mussel, *Mytilus edulis* is part of the *Mytilus edulis* species complex, encompassing at least three putative species: *M. edulis, M. galloprovincialis* and *M. trossulus*. These three species occur on both sides of the Atlantic and hybridize in nature, and both *M. edulis* and *M. galloprovincialis* are important aquaculture species. They are also invasive species in many parts of the world. Here, we present a chromosome-level assembly of *Mytilus edulis* . We used a combination of PacBio sequencing and Dovetail's Omni-C technology to generate an assembly with 14 long scaffolds containing 94% of the predicted length of the *M. edulis* genome (1.6 out of 1.7 Gb). Assembly statistics were total length 1.65 Gb, N50 = 116 Mb, L50 = 7 and, L90 = 13. BUSCO analysis showed 92.55% eukaryote BUSCOs identified. AB-*Initio* annotation using RNA-seq from mantle, gills, muscle and foot predicted 47,128 genes. These gene models were combined with Isoseq validation resulting in 65,505 gene models and 129,708 isoforms. Using GBS and shotgun sequencing, we also sequenced 3 North American populations of *Mytilus* to characterize single-nucleotide as well as structural variance. This high-quality genome for *M. edulis* provides a platform to develop tools that can be used in breeding, molecular ecology and evolution to address questions of both commercial and environmental perspectives.
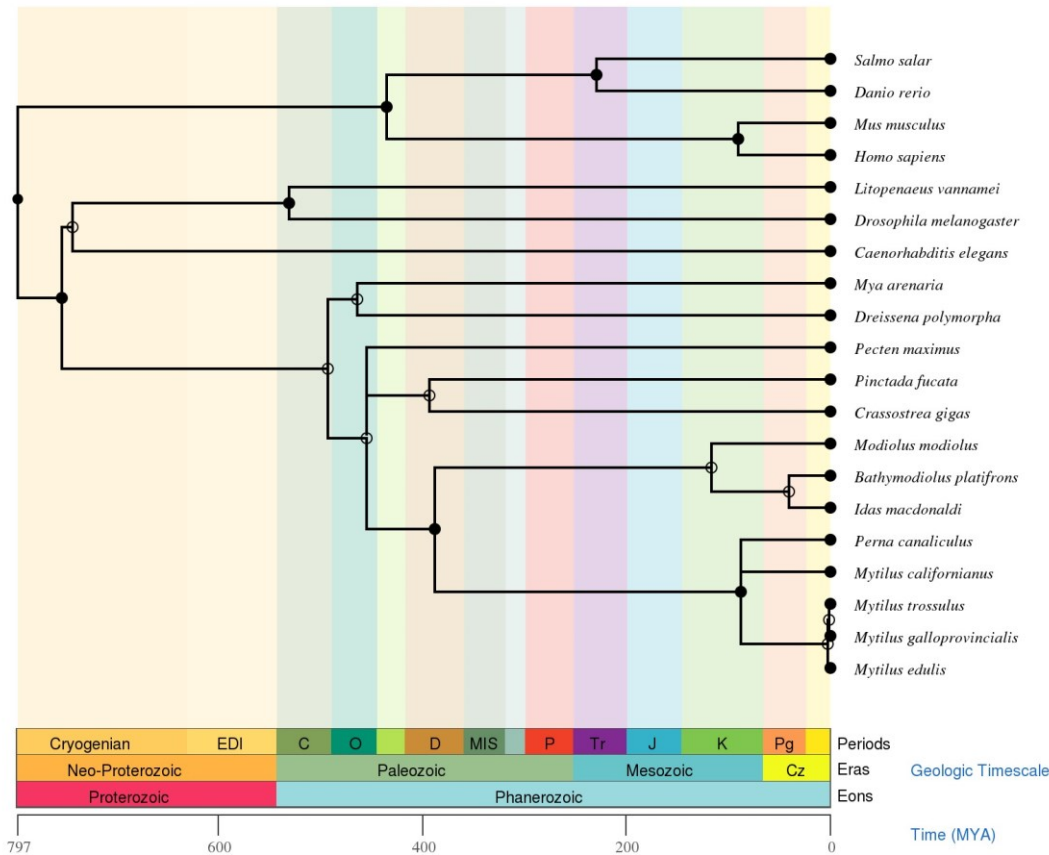
26    **Introduction**

27    The blue mussel (*Mytilus edulis*) is common to the North Atlantic from Arctic to Mediterranean

28    regions, with habitat ranging from upper shore to the shallow subtidal (Hayward and Ryland 2017).

29    *M. edulis* is known to hybridise with *M. trossulus* in North America and with *M. trossulus* and *M.*

30    *galloprovincialis* in Europe. Together, these three species form the *M. edulis* species complex (**Fig. 1**)

31    (McDonald et al. 1991).

32    This reef-building bivalve is an ecosystem engineer. Blue mussels dominate fouling communities in

33    shallow and substrata providing important secondary habitat (Norling and Kautsky 2007). Offshore

34    wind energy structure surveys found that they can cover the structures with up to 3.4 kg of biomass

35    m$^{-2}$ (Krone et al. 2013). Through filter-feeding, eutrophication is reduced which can alter ecosystems

36    (Broszeit et al. 2016). This nutrient cycling ability has been harnessed by using *M. edulis* to study the

37    fate of persistent organic and metal pollutants (Chase et al. 2001; McEneff et al. 2014), for the

38    bioremediation of waste (Broszeit et al. 2016) and reduce environmental effects from salmon farms

39    (MacDonald et al. 2011).

40    Mussels, a key bivalve production species (FAO 2020), face decreasing wild spat availability for

41    aquaculture in the UK and elsewhere (Regan et al. 2021). These losses are attributed to multiple

42    stressors including warmer seas (Seuront et al. 2019) causing a poleward range contraction (Jones et

43    al. 2010). Additionally, warmer oceans elevate dissolved $CO_2$ leading to Ocean Acidification (OA)

44    impacting mussel viability (Asplund et al. 2014) and disease resistance (Ellis et al. 2015). This makes

45    mussels more susceptible to bacterial pathogenesis (Eggermont et al. 2017; Ripabelli et al. 1999), with

46    emerging pathogens posing a constant threat (Charles et al. 2020; Cano et al. 2022). Infectious disease

47    such as disseminated Neoplasia (DN) of *M. trossulus* origin is associated with reduced fitness in *Mytilus*

48    *spp.* (Burioli et al. 2021). Furthermore, the effects of hybridisation between the three species of the

49    *Mytilus edulis* complex are yet uncertain with suggested negative effect of *M. trossulus* hybridisation

50    and potential adaptive introgression in the case of *M. galloprovincialis* hybridisation (Fraïsse et al.

51    2014; Kenchington et al. 2020; Michalek et al. 2021).

52    To protect the aquaculture industry from these threats, hatchery efforts have been launched in the

53    UK (Regan et al. 2021), elsewhere in Europe (Kamermans et al. 2013) and in Canada (Gurney-Smith et

54    al. 2017). However, these efforts have not been straightforward and a better understanding of

55    fundamental biology is required to achieve commercial success. Despite their importance in

56    aquaculture and the valuable ecosystem services they provide, no chromosomal assembly existed for

57    any species within the *Mytilus edulis* species complex prior to this study. Improved genomic tools are

58    required to address fundamental biological questions such as inheritance patterns and adaptations.

59    Like many bivalves, the mussel genome is highly heterozygous (3.5%) with an estimated 43% repeat

60    content. The linear plot of k-mer abundance analysis clearly shows a heterozygous peak in addition to

61    the homozygous peak and estimates a shorter haploid genome length of 1.18 Gb compared to flow-

62    cytometry data. The estimated repeat content of the *M. edulis* genome is ~43%. These characteristics

63    make assembly of these genomes challenging. However, recently genomes for the American oyster,

64    Pacific Oyster, *M. corruscus* and *M. califoniaus* have been assembled to chromosome-level using short

65    and long-read sequencing technologies as well as Hi-C-based scaffolding (Yang et al. 2021; Paggeot et

66    al. 2022a; Penaloza et al. 2021; Gomez-Chiarri et al. 2015). We used a similar approach in this project

67    to produce a highly contiguous assembly. Practical application of this assembly is demonstrated in

68    cross-species synteny analyses and in population structure of *Mytilus* individuals sampled from

69    different regions of the Canadian Atlantic.

**Fig. 1. TimeTree for *Mytilus edulis***

Mytilidae diverged ~387 MYA. Generated using TimeTree (Kumar et al. 2022).

74

## Methods

*DNA extraction, library preparation for genome assembly*

One naïve blue mussel sample (Anne) was selected from samples collected by the Provincial Department of Communities and Fisheries in the estuary Foxley river in PE. Foxley River was selected as a sampling site because there is no grow-out aquaculture there (i.e. seed from other bays are not transferred to the area). Due to potential introgression of other species of the *Mytlius* species complex (e.g. *M. trossulus*), this sample was genotyped using 12 SNPs described by (Wilson et al. 2018). We sampled the gill, mantle and muscle of sample "Anne" aseptically, flash froze fresh tissues in liquid nitrogen and preserved them at -80°C. Tissues were shipped to Dovetail Genomics in Scott's Valley, CA, in excess dry ice. Dovetail extracted high molecular weight (HMW) DNA using an in-house modified CTAB method. Dovetail prepared and sequenced PacBIO SMRTbell libraries to a depth of 196X using the Sequel II sequencer. Processed PacBio data (from Dovetail) can be found on SRA under accession number SRX11246493.

*Raw contig assembly, scaffold formation and polishing*

Dovetail generated a primary contig-level assembly using wtdbg2 (Ruan and Li 2020). The contig-level assembly was filtered of putative duplicated haplotypes and contaminants using Purge Haplotigs (Roach et al. 2018) and Blobtools2 (Challis et al. 2020), respectively. Scaffolding was performed using Omni-C libraries and the HiRise assembler (v1.0). The same DNA sample used by Dovetail was shipped to UWM (University of Wisconsin – Madison) and sequenced to a 100X using the NovaSeq sequencer. These data were trimmed using Trimmomatic (v) and used for polishing with racon (v1.4.3) (Vaser et al. 2017). Completeness was evaluated using compleasm v0.2.5 (Huang and Li 2023) (eukaryota_odb10: 255 BUSCOS, metazoa_odb10: 954 BUSCOS and mollusca_odb10: 5295 BUSCOS) (Manni et al. 2021) and Merqury (Rhie et al. 2020). Merqury analysis was carried out with the same read set used for polishing as the original PacBio CLR reads were not suitable for this analysis. General

99    length metrics were obtained using QUAST (v5.0.1) (Gurevich et al. 2013; Mikheenko et al. 2018).

100    Synteny mapping between the *M. edulis* and the *M. coruscus* (GCA_017311375.1) (Yang et al. 2021)

101    was done using the MCScanX.h function of MCScanX (v2) (Wang et al. 2012). Putative orthologous

102    groups were identified with Orthofinder (v.2.5.4) using predicted gene structures for *M. edulis* (this

103    work) and *M. coruscus* (Yang et al. 2021). Dot plots and circle plots were generated using MCScanX

104    (v.2).

105    *RNA preparation and Isoseq analysis*

106    Isoseq3 analysis of CSS data from muscle, gill, hemolymph, and foot (2.9 million reads -188,165 Mb)

107    identified 216,343 high-quality putative full-length (FL) transcripts (from 2.8 million reads containing

108    poly-A tails). We shipped flash frozen gill and adductor muscle tissue from sample "Anne" to the

109    Biotechnology Centre Core facility at the University of Wisconsin, Madison (UWM). Samples were

110    homogenized using a Qiagen Tissuelyser (2 min @ 20 Hz). RNA was extracted using the RNeasy Mini

111    Kit (Qiagen) with on-column DNAse treatment. UWM performed RNA QC with a nanodrop and

112    bioanalyzer and prepared libraries using the Iso-seq Library SMRTbell express template prep kit.

113    Libraries were sequenced using one Sequel II SMRT cell in CSS mode (i.e. HiFi reads). UWM provided

114    de-multiplexed processed HiFi reads. RNA samples from foot, mantle and gut were sent to the

115    Genome Excellence Centre (Genome Quebec) in Montreal. HiFi sequencing was carried out as above.

116    Putative full-length transcripts were identified using the IsoSeq3 pipeline

117    (https://github.com/ylipacbio/IsoSeq3). Putative open read frames (ORF) were identified using

118    TransDecoder (v5.5.0) (https://github.com/TransDecoder/TransDecoder/wiki).

119    *Annotation*

120    Repeat modeller (4.1.0) (Flynn et al. 2020) was used to predict repeat motifs for *M. edulis* and Repeat

121    Masker (2.0.2a) (Smit et al. 2015) was used to mask the final assembly. Ab-initio annotation was done

122    using Augustus (v3.4.0) (Stanke et al. 2006) trained with genes from *Mytilus galloprovincialis*, *Mytilus*

123    *coruscus* and *Crassostrea virginica*. Additional hints were generated using short-read RNA-seq (data

124    not shown) and Isoseq data generated herein. The ab-initio annotation was updated using PASA (v2.5)

125    with alignments of a de-novo transcriptome produced using Trinity (v2.8.15) (Grabherr et al. 2011) in

126    Genome-Guided mode. Two runs through PASA were used to update the *ab-initio* annotation. Full-

127    length (FL) transcripts from Isoseq3 pipeline were mapped to the genome using pbmm2 with pre-sets

128    for Isoseq and filtered based on quality Isoseq3 collapse (minimum alignment identity/coverage:

129    0.90/0.90) and Isoseq3 refine. We used the resulting gff annotation to run SQANTI3 (v4.3.0). We used

130    sqanti_qc.py to generate quality data for sqanti_filter.py. Filtering to remove artifacts was carried out

131    using the default parameters. Lastly, ab-initio predications and filtered Isoseq FLs were merged with

132    AGAT (v.1.2.0). Using agat_sp_merge.pl, we removed duplicate gene models/isoform and assigned

133    orphan isoforms from the Isoseq data when possible. Amino acid sequences were translated from the

134    CDS using agat_sp_extract_sequences using options -t "CDS" -p --cfs --acs –asc.

135    *Sample collection for SNP discovery and population structure*

136    We collected gill samples for DNA extractions using standard molecular biology techniques and

137    preserved them in non-denatured 100% ethanol. For population genetics analysis, we sampled

138    mussels in sets of 96 from west to east PEI in Foxley River (FOX) (wild population), Malpeque Bay (MB)

139    (North Shore, Prince County); French River (FR), Stanley Bridge (ST), Whetley River (WR), Tracadie

140    (TRC) (North Shore, Queen county); Orwell (OR) (South Shore, Queen's county; Morell (MRL) (North

141    Shore – King's county) and; Murray River (MR) (South Shore, King county) (**Fig. 2**). Seed deployed on

142    these sites were originally collected in St. Peter's bay, Brudenell River, and Malpeque bay. We also

143    collected samples in the Bras D'or Lake in Cape Breton (Nova Scotia), the Magdalene Island (Québec)

144    and Notre Dame bay in Newfoundland.

145    *SNP markers, Admixture analysis, population structures.*

146    Samples for SNP discovery and population genetics were shipped to LGC genomics in Berlin.

147    Restriction-associated DNA libraries (RAD) libraries were prepared by LGC with using MslI, normalized

148    and sequenced using a NextSeq sequencer. The resulting reads were trimmed and checked for the

149    restriction site by LGC. This final read set was used to identify SNPs and call individual genotypes using

150    Tassel5 (v.5.2.4). SNPs were filtered using MAF (<0.01) and percent of individuals with genotypes. For

151    the SNP discovery, the samples from Cape Breton were excluded from filtering analysis. These samples

152    were shown to be a pure *M. trossulus* population and had missing calls for a significant number of

153    sites. For population genetics analysis, a second set of SNPs that were successfully called across all

154    populations was used. All samples were also genotyped for the 12 species discrimination SNPs from

155    (Wilson et al. 2018). All population genetics analyses were performed using the dpca function of the

156    R package Adegenet (v.2.1.7) and STRUCTURE (v.2.3.4) or fastSTRUCTURE (v1.0) (K 3 to K10, 5000

157    repetitions, 1000 burn in). Analysis of species discrimination SNPs included the genotypes published

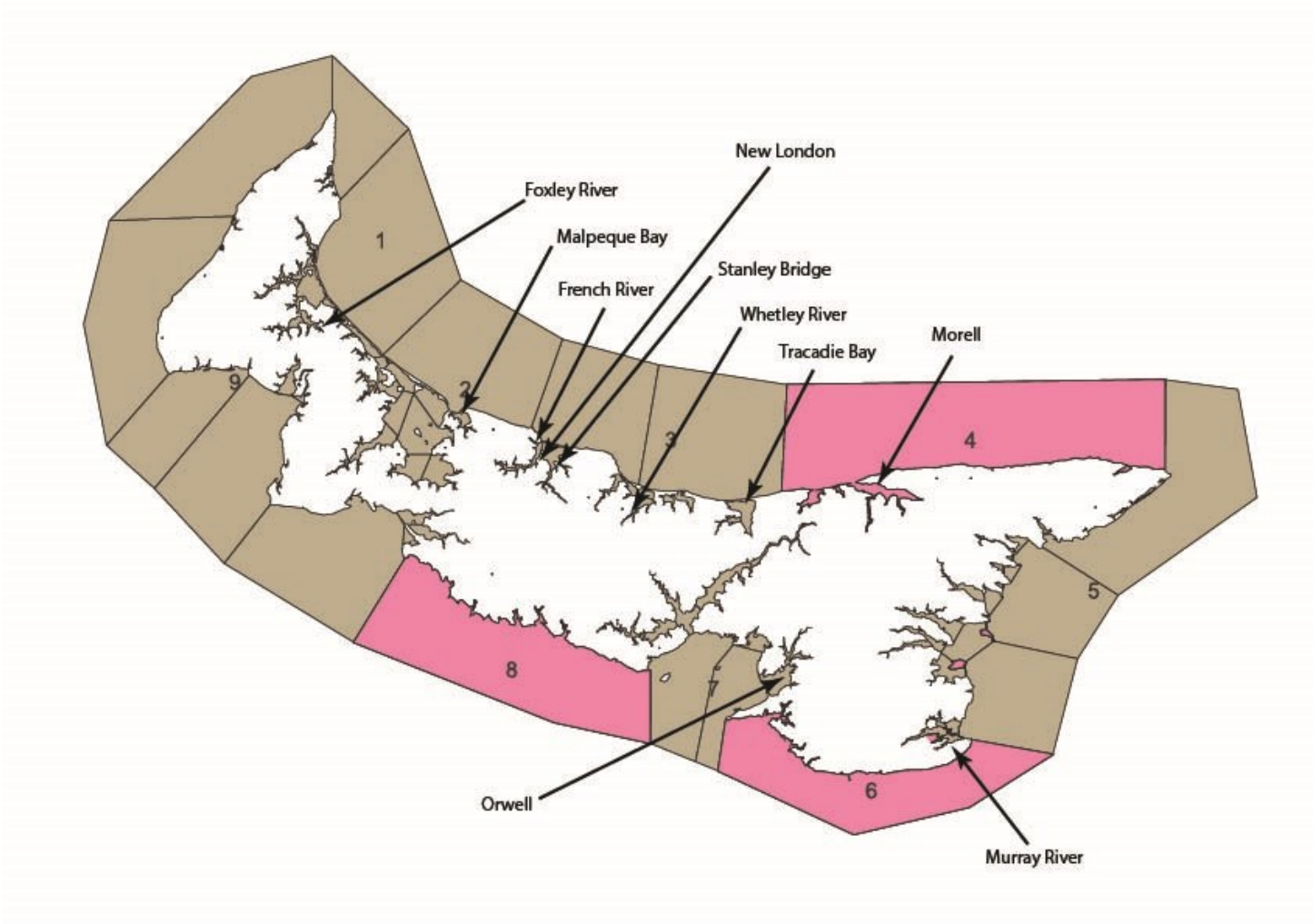158    by (Corrochano-Fraile et al. 2022) as outgroups.

159

160



Figure 2: Sampling Locations in PEI

161

162

163

164 **Results and Discussion**

165 *Genome Assembly and Annotation*

166 The chromosome-level assembly presented herein was produced in two stages. First, ~15 million

167 PacBio CLR reads (~340 Gb) were produced, representing coverage of 196X for an estimated genome

168 size of 1.7 Gb (Hinegardner 1974; Rodríguez-Juíz et al. 1996). These reads were assembled into contigs

169 using wtdbg2, which uses uncorrected reads (Ruan and Li 2020). The primary assembly was 1.96 Gb

170 long in 17,825 contigs and a N50 of 443 Kb. After haplotype purging and contaminant removal, the

171 final contig assembly had 10,111 contigs, for a total of 1,65 Gb and a N50 of 518 Kb. Following

172 scaffolding using Omni-C libraries and the HiRise assembler, we generated a primary chromosome-

173 level assembly made of 2,117 contigs. We removed putative contaminants using Blobtools by

174 eliminating sequences coming from non-molluscan organisms. This assembly was further filtered to

175 contain only sequences > 5,000 bp. The resulting draft is deposited on NCBI assembly under accession

176 number GCA_019925275.1. The final assembly is made of 1,119 contigs and has an N50 of 116 Mb.

177 The 14 putative *M. edulis* chromosomes are deposited under accession numbers CM034349.1 to

178 CM34362.1. Detailed statistics for the assemblies can be found in **Table 1**.

179

| | Anne (This study) | Corrochano-Fraile *et al.* (2022) | xbMytEdul2.1, Darwin Tree of Life (2024) |
|---|---|---|---|
| *Length* | 1,659,567,081 bp | 1,827,085,763 | 1,374,471,240 |
| *# of scaffolds* | 1,119 | 3,339 | 2,563 |
| *# of contigs* | 9,866 | 5,966 | 3,754 |
| *N50* | 116,503,180 | 1,097,279 | 1,734,586 |
| *Coverage* | 196.0x | 152.0x | 30.0x |
| *Completeness (Merqury)* | 76.71 | NA | NA |
| *QV* | 32.39 | NA | NA |

180 **Table 1. Summary statistics for *M. edulis* assemblies**

181 Compared to the assembly published in Corrochano-Fraile et al. (2022), the assembly presented

182 herein has better contiguity (**Table 1**). Our assembly is also shorter than the assembly from

183   Corrochano-Fraile et al. (2022). The two assemblies' length falls close to the estimated size of the

184   genome based on c-value (1.7) (Rodríguez-Juíz et al. 1996) and is significantly longer than what is

185   estimated by k-mer abundance analysis with GenomeScope (1.18 Gb).

186   Despite the total assembly length being close to that estimated using c-values, k-mer-based

187   completeness analysis recovers only 76% in a set of Illumina reads origination from sample "Anne".

188   When the putative purged haplotigs were added back to the assembly, recovery was ~83%. This

189   apparent low k-mer recovery is probably a combination of the consensus being different from either

190   haploid genomes, the error rate in the original PacBio data, the fact that the reads used for polishing

191   were not used from the primary assembly, the contigs removed based on length or contaminant

192   status, and the gaps arising from the Omni-C scaffolding. It is also possible that the high heterozygosity

193   affects the accuracy of k-mer abundance analysis, as shown by the large discrepancy between genome
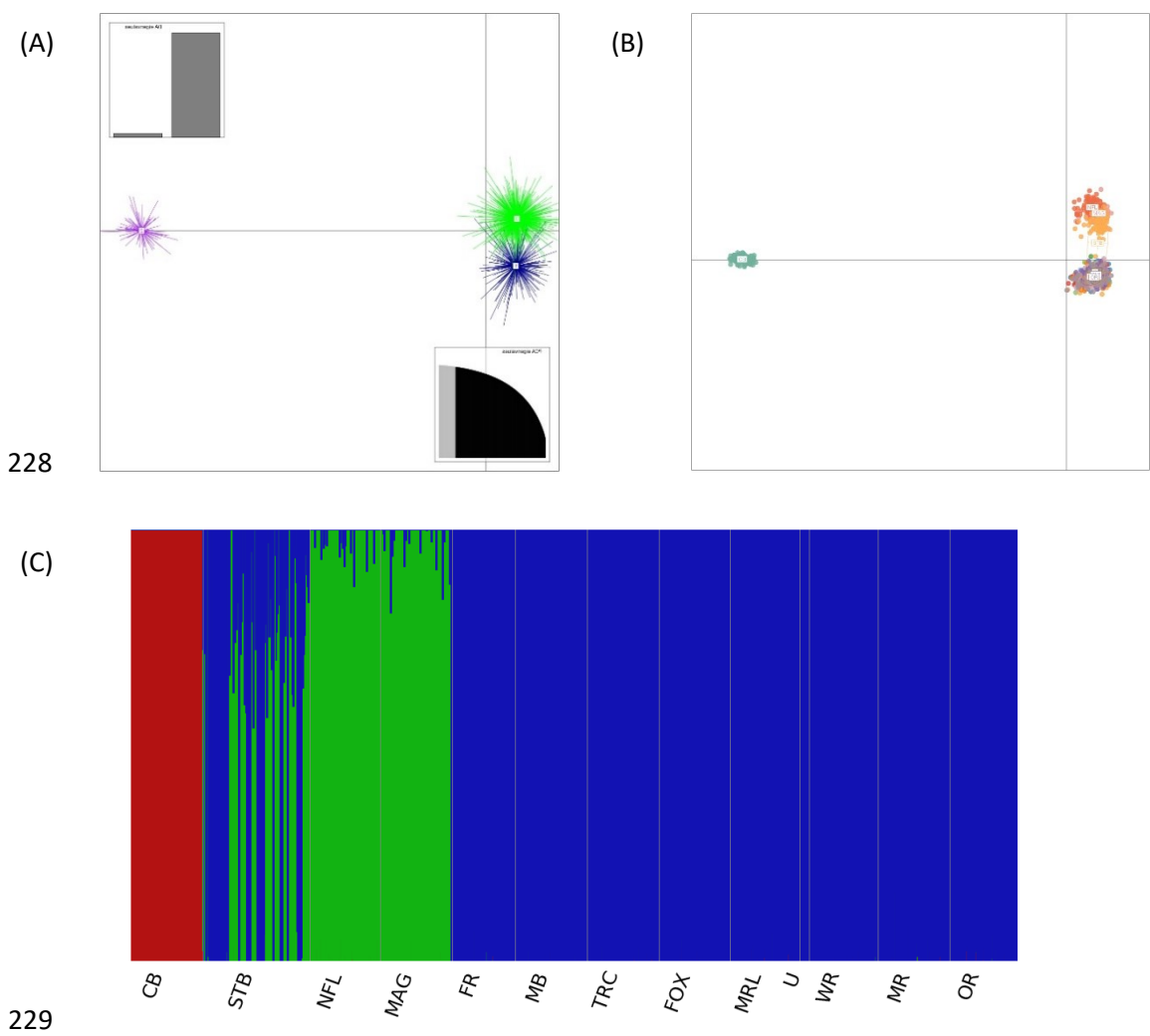
194   size estimates.

195   Completeness analysis in Merqury resulted in 76.71% recovery of k-mers present in the polishing

196   Illumina data from the final version of the assembly. We also evaluated the completeness of the

197   assembly when combined with purged haplotypes, which was 83.44%. QV value for the primary

198   assembly was 32.39, while the combined draft had a QV of 30.73. Compleasm BUSCO analysis showed

199   a recovery of 92.55%, 91.95%, and 88.5% complete BUSCOS against the eukaryote, metazoan and

200   molluscan databases, respectively. The 14 putative chromosomes represent ~96% of the assembly,

201   with lengths varying from 140 Mbp to 90 Mbp. These data and the N50 metric show that this assembly

202   has high contiguity and that this assembly and its annotation will be highly useful for aquaculture,

203   evolution and molecular ecology studies. Herein, we illustrate the possible applications of this

204   assembly by performing population and synteny analyses.

205   *SNP discovery and population structure*

206   Due to the close relationship between the members of the *Mytilus* species complex, we wanted to

207   verify that the individual sampled (Anne) was pure *M. edulis*. Population Structure analysis conducted

208    using 12 SNPs (Wilson et al. 2018) clearly separated *M. edulis* populations in PEI from other regions of

209    Canada. We generated two sets of SNPs: the first set totalling 71,231 SNPs using only samples from

210    PEI and made a polymorphic collection of SNPs in *M. edulis.* The second set, with ~6,000 markers, is a

211    polymorphic set of SNPs in both *M. edulis* and *M. trossulus.* Population structure and putative

212    admixture are shown in **Figure 4**. In the DPCA and PCA, untrained clustering clearly separated both CB

213    from PEI/NL/MAG and also the putative populations in the Gulf of St. Lawrence. In green, are the

214    majority of samples from NL and MAG, while blue represents individuals from PEI.

215    We used *M. trossulus*, *M. galloprovincialis* and European *M. edulis* genotypes as outgroups. The later

216    inclusion of 96 individuals from Cape Breton (NS) showed no significant evidence of *M. trossulus*

217    introgression in PEI samples. Given that Cape Breton has long been considered a pure *M. trossulus*

218    population (Wilson et al., 2018), we are confident that the sample Anne represents an *M. edulis*

219    individual. We also genotyped over 500 PEI individuals and 96 samples from Magdalene Island and 96

220    individuals from Newfoundland with the same 12 SNP panel. As before no significant introgression of

221    *M. trossulus* was detected in PEI. However, we only genotyped 96 samples collected in an area with

222    no grow-out leases (Foxley River). Although unlikely, we cannot rule out the possibility of a sampling

223    bias in aquaculture sites favouring *M. edulis.* DPCA and Structure analysis indicate that there is low

224    population stratification between different regions of PEI, while the Magdalene islands and NL are

225    distinct populations. The populations from NL and the Magdalene Islands are more similar to each

226    other than from the populations from PEI. For animals from one sampling event at S.t Peter's bay, we

227    found evidence of shared genetics between PEI and the population to the northeast of the island.

(A)

(B)

228

(C)

229

**Figure 4 Population structure in the North American Atlantic.** (A) DPCA with 6,000 SNPs, (B) PCA with 6,000 SNPs, (C) fastStructure with 6,000 SNPs. Samples from Cape Breton NS (CB), FR, TRC, MB, FOX, MRL, U, WR, MR, OR (PEI), Magdalene Islands QC (MAG), Notre Dame Bay NL (NFL)

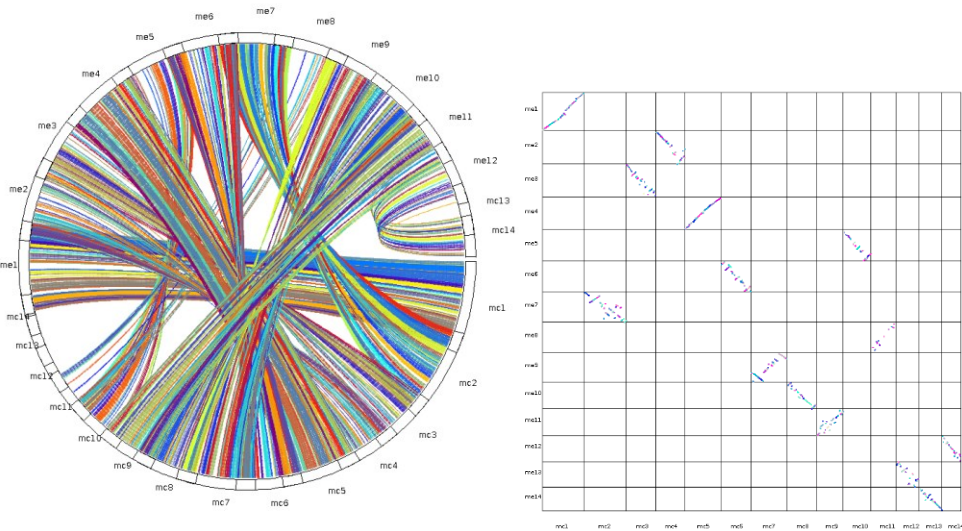234    *Annotation and synteny analysis*

235    We identified 196,111 putative open reading frames from the 216,343 FL transcripts using Isoseq3

236    analysis of CSS data from muscle and gill (2.9 million reads -188,165 Mb). BLASTp analysis against the

237    uniref90 database returned informative hits for ~80% (164,969) of these translated transcripts. *Ab-*

238    *initio* gene prediction in Augustus detected 46,604 gene models that produced 46,604 transcripts after

239    filtering based on evidence support. After two rounds of PASA updates, the final number of gene

240    models was 47,128 and the number of transcripts was 55,138. After Isoseq refine and collapse 85,099

241    isoforms survived and following SQANTI3 filtering 70,592 isoforms remained. They were assigned to

242    31,211 unique genes. Finally, the combined annotation has 65,505 gene models and 129,708 isoforms.

243    Proteins were translated from the CDS ensuring only complete CDS were translated and that isoforms

244    were not incorrectly fused together. This resulted in 45,379 amino acid sequences. Compleasm BUSCO

245    analysis of these 45,379 proteins showed a recovery of 78.43% and 76.83% complete BUSCOS against

246    the eukaryote and metazoan databases, respectively.

247    Synteny analysis showed a high degree of collinearity between putative chromosomes of *M. edulis*

248    and *M. coruscus* (**Fig. 3**). However, putative inversions, transposition and deletion can be observed in

249    almost all chromosomes. Gene order in chromosomes represented by sequences CM034349.1 (*M.*

250    *edulis* chromosome 1) - CM029595.1 (*M. coruscus* LG1) and CM034343.1 (*M. edulis* chromosome 4) –

251    CM029599.1 (*M. coruscus* LG5) showed the highest conservation. Putative orthologous relationships

252    between *M. edulis* chromosomes and *M. coruscus* linkage groups (LG) are shown in **Table 2**.

253    The taxonomic status of the "species" in the *Mytilus* species complex remains in debate. Chromosome-

254    level assemblies allow the study of macroevolution of the genome by looking at synteny across

255    species. Herein, we present the synteny analysis between the 14 putative chromosomes of *M. edulis*

256    and *M. corruscus* (Yang et al. 2021) to exemplify how chromosome-level assemblies may allow us to

257    better understand the phylogenetic relationships within the genus *Mytilus*. This analysis shows that

258    some of putative orthologous chromosomes of the 2 species maintain high-levels of collinearity (e.g.

259     chromosomes 1 and 4 from *M. edulis* with LG1 and LG5 from *M. coruscus* respectively) while others

260     present significant levels of re-arrangements (e.g. chromosomes 7 and 11 from *M. edulis* with LG2 and

261     LG9 from *M. coruscus* respectively). An in-depth analysis of chromosome synteny will shed light on

262     the level of collinearity between multiple members of the genus *Mytilus*. Homology between

263     chromosomes is a key element of the viability of hybrids. Reproductive isolation tends to increase

264     during speciation, and these resources will permit further studies on the reproductive compatibility

265     of the species in genus *Mytilus* at the chromosome level.

266

267



268

**Figure 3. Synteny between the *M. edulis* and *M. coruscus* putative chromosomes**

270

**Table 2: Putative synteny between *M. edulis* and *M. coruscus.* Ids are NCBI Assembly database Molecule name and accession number.**

| *M. edulis* | *M. coruscus* |
|---|---|
| Chromosome 1 (CM034349.1) | LG01 (CM029595.1) |
| Chromosome 2 (CM034350.1) | LG04 (CM029598.1) |
| Chromosome 3 (CM034351.1) | LG03 (CM029597.1) |
| Chromosome 4 (CM034352.1) | LG05 (CM029599.1) |
| Chromosome 5 (CM034353.1) | LG10 (CM029604.1) |
| Chromosome 6 (CM034354.1) | LG06 (CM029600.1) |
| Chromosome 7 (CM034355.1) | LG02 (CM029596.1) |
| Chromosome 8 (CM034356.1) | LG11 (CM029605.1) |
| Chromosome 9 (CM034357.1) | LG07 (CM029601.1) |
| Chromosome 10 (CM034358.1) | LG08 (CM029602.1) |
| Chromosome 11 (CM034359.1) | LG09 (CM029603.1) |
| Chromosome 12 (CM034360.1) | LG14 (CM029608.1) |
| Chromosome 13 (CM034361.1) | LG12 (CM029606.1) |
| Chromosome 14 (CM034362.1) | LG13 (CM029607.1) |

273

274

275    Here, we present a highly contiguous chromosome assembly for *Mytilus edulis* confirming species-

276    level individual purity through resequencing. To date, our resource has been applied in multiple

277    studies analysing *Mytilus* genome assemblies (Paggeot et al. 2022b; Gallardo-Escarate et al. 2023) and

278    cross-species gene orthology analyses (Saco et al. 2023). The gene annotations produced in this study

279    were generated using Augustus gene model predictions integrating full transcript Isoseq data and

280    applying stringent filtering parameters. This comprehensive approach provides a robust foundation

281    for future cross-species analyses and biological studies on gene function within the *Mytilus* species

282    complex.

283

300 Asplund, M.E., S.P. Baden, S. Russ, R.P. Ellis, N. Gong *et al.*, 2014 Ocean acidification and host-
301     pathogen interactions: blue mussels, Mytilus edulis, encountering Vibrio tubiashii. *Environ*
302     *Microbiol* 16 (4):1029-1039.
303 Broszeit, S., C. Hattam, and N. Beaumont, 2016 Bioremediation of waste under ocean acidification:
304     Reviewing the role of Mytilus edulis. *Marine Pollution Bulletin* 103 (1):5-14.
305 Burioli, E.A.V., M. Hammel, N. Bierne, F. Thomas, M. Houssin *et al.*, 2021 Traits of a mussel
306     transmissible cancer are reminiscent of a parasitic life style. *Scientific Reports* 11 (1):24110.
307 Cano, I., A. Parker, G.M. Ward, M. Green, S. Ross *et al.*, 2022 First Detection of Francisella halioticida
308     Infecting a Wild Population of Blue Mussels Mytilus edulis in the United Kingdom. *Pathogens*
309     11 (3).
310 Challis, R., E. Richards, J. Rajan, G. Cochrane, and M. Blaxter, 2020 BlobToolKit – Interactive Quality
311     Assessment of Genome Assemblies. *G3 Genes|Genomes|Genetics* 10 (4):1361-1374.
312 Charles, M., A. Villalba, G. Meyer, S. Trancart, C. Lagy *et al.*, 2020 First detection of Francisella
313     halioticida in mussels Mytilus spp. experiencing mortalities in France. *Dis Aquat Organ*
314     140:203-208.
315 Chase, M., S.H. Jones, P. Hennigar, J. Sowles, G. Harding *et al.*, 2001 Gulfwatch: Monitoring spatial
316     and temporal patterns of trace metal and organic contaminants in the Gulf of Maine (1991–
317     1997) with the blue mussel, Mytilus edulis L. *Marine Pollution Bulletin* 42 (6):490-504.
318 Corrochano-Fraile, A., A. Davie, S. Carboni, and M. Bekaert, 2022 Evidence of multiple genome
319     duplication events in Mytilus evolution. *BMC Genomics* 23 (1):340.
320 Eggermont, M., P. Bossier, G.S.J. Pande, V. Delahaut, A.M. Rayhan *et al.*, 2017 Isolation of
321     Vibrionaceae from wild blue mussel (Mytilus edulis) adults and their impact on blue mussel
322     larviculture. *FEMS Microbiology Ecology* 93 (4).
323 Ellis, R.P., S. Widdicombe, H. Parry, T.H. Hutchinson, and J.I. Spicer, 2015 Pathogenic challenge
324     reveals immune trade-off in mussels exposed to reduced seawater pH and increased
325     temperature. *Journal of Experimental Marine Biology and Ecology* 462:83-89.

326  FAO, 2020 The State of World Fisheries and Aquaculture 2020. Sustainability in action. , pp. #244 p.
327      in *The State of World Fisheries and Aquaculture (SOFIA)*. FAO, Rome, Italy.
328  Flynn, J.M., R. Hubley, C. Goubert, J. Rosen, A.G. Clark *et al.*, 2020 RepeatModeler2 for automated
329      genomic discovery of transposable element families. *Proceedings of the National Academy*
330      *of Sciences* 117 (17):9451-9457.
331  Fraïsse, C., C. Roux, J.J. Welch, and N. Bierne, 2014 Gene-flow in a mosaic hybrid zone: is local
332      introgression adaptive? *Genetics* 197 (3):939-951.
333  Gallardo-Escarate, C., V. Valenzuela-Munoz, G. Nunez-Acuna, D. Valenzuela-Miranda, F.J. Tapia *et al.*,
334      2023 Chromosome-Level Genome Assembly of the Blue Mussel Mytilus chilensis Reveals
335      Molecular Signatures Facing the Marine Environment. *Genes (Basel)* 14 (4).
336  Gomez-Chiarri, M., W.C. Warren, X. Guo, and D. Proestou, 2015 Developing tools for the study of
337      molluscan immunity: The sequencing of the genome of the eastern oyster, Crassostrea
338      virginica. *Fish Shellfish Immunol* 46 (1):2-4.
339  Grabherr, M.G., B.J. Haas, M. Yassour, J.Z. Levin, D.A. Thompson *et al.*, 2011 Full-length
340      transcriptome assembly from RNA-Seq data without a reference genome. *Nature*
341      *Biotechnology* 29 (7):644-652.
342  Gurevich, A., V. Saveliev, N. Vyahhi, and G. Tesler, 2013 QUAST: quality assessment tool for genome
343      assemblies. *Bioinformatics* 29 (8):1072-1075.
344  Gurney-Smith, H.J., A.J. Wade, and C.L. Abbott, 2017 Species composition and genetic diversity of
345      farmed mussels in British Columbia, Canada. *Aquaculture* 466:33-40.
346  Hayward, P.J., and J.S. Ryland, 2017 *Handbook of the marine fauna of North-West Europe*: Oxford
347      University Press.
348  Hinegardner, R., 1974 Cellular DNA content of the Mollusca. *Comp Biochem Physiol A Comp Physiol*
349      47 (2):447-460.
350  Huang, N., and H. Li, 2023 compleasm: a faster and more accurate reimplementation of BUSCO.
351      *Bioinformatics* 39 (10).
352  Jones, S.J., F.P. Lima, and D.S. Wethey, 2010 Rising environmental temperatures and biogeography:
353      poleward range contraction of the blue mussel, Mytilus edulis L., in the western Atlantic.
354      *Journal of Biogeography* 37 (12):2243-2259.
355  Kamermans, P., T. Galley, P. Boudry, J. Fuentes, H. McCombie *et al.*, 2013 11 - Blue mussel hatchery
356      technology in Europe, pp. 339-373 in *Advances in Aquaculture Hatchery Technology*, edited
357      by G. Allan and G. Burnell. Woodhead Publishing.
358  Kenchington, E.L., B.W. MacDonald, A. Cogswell, L.C. Hamilton, and A.P. Diz, 2020 Sex-specific
359      effects of hybridization on reproductive fitness in Mytilus. *Journal of Zoological Systematics*
360      *and Evolutionary Research* 58 (2):581-597.
361  Krone, R., L. Gutow, T.J. Joschko, and A. Schröder, 2013 Epifauna dynamics at an offshore foundation
362      – Implications of future wind power farming in the North Sea. *Marine Environmental*
363      *Research* 85:1-12.
364  Kumar, S., M. Suleski, J.M. Craig, A.E. Kasprowicz, M. Sanderford *et al.*, 2022 TimeTree 5: An
365      Expanded Resource for Species Divergence Times. *Mol Biol Evol*.
366  MacDonald, B.A., S.M.C. Robinson, and K.A. Barrington, 2011 Feeding activity of mussels (Mytilus
367      edulis) held in the field at an integrated multi-trophic aquaculture (IMTA) site (Salmo salar)
368      and exposed to fish food in the laboratory. *Aquaculture* 314 (1):244-251.
369  Manni, M., M.R. Berkeley, M. Seppey, and E.M. Zdobnov, 2021 BUSCO: Assessing Genomic Data
370      Quality and Beyond. *Current Protocols* 1 (12):e323.
371  McDonald, J., R. Seed, and R. Koehn, 1991 Allozymes and morphometric characters of three species
372      ofMytilus in the Northern and Southern Hemispheres. *Marine Biology* 111 (3):323-333.
373  McEneff, G., L. Barron, B. Kelleher, B. Paull, and B. Quinn, 2014 A year-long study of the spatial
374      occurrence and relative distribution of pharmaceutical residues in sewage effluent, receiving
375      marine waters and marine bivalves. *Science of the Total Environment* 476:317-326.

376     Michalek, K., D.L.J. Vendrami, M. Bekaert, D.H. Green, K.S. Last *et al.*, 2021 Mytilus trossulus
377          introgression and consequences for shell traits in longline cultivated mussels. *Evolutionary*
378          *applications* 14 (7):1830-1843.
379     Mikheenko, A., A. Prjibelski, V. Saveliev, D. Antipov, and A. Gurevich, 2018 Versatile genome
380          assembly evaluation with QUAST-LG. *Bioinformatics* 34 (13):i142-i150.
381     Norling, P., and N. Kautsky, 2007 Structural and functional effects of Mytilus edulis on diversity of
382          associated species and ecosystem functioning. *Marine Ecology Progress Series* 351:163-175.
383     Paggeot, L.X., M.B. DeBiasse, M. Escalona, C. Fairbairn, M.P.A. Marimuthu *et al.*, 2022a Reference
384          genome for the California ribbed mussel, Mytilus californianus, an ecosystem engineer. *J*
385          *Hered*.
386     Paggeot, L.X., M.B. DeBiasse, M. Escalona, C. Fairbairn, M.P.A. Marimuthu *et al.*, 2022b Reference
387          genome for the California ribbed mussel, Mytilus californianus, an ecosystem engineer. *J*
388          *Hered* 113 (6):681-688.
389     Penaloza, C., A.P. Gutierrez, L. Eory, S. Wang, X. Guo *et al.*, 2021 A chromosome-level genome
390          assembly for the Pacific oyster Crassostrea gigas. *Gigascience* 10 (3).
391     Regan, T., T.P. Bean, T. Ellis, A. Davie, S. Carboni *et al.*, 2021 Genetic improvement technologies to
392          support the sustainable growth of UK aquaculture. *Reviews in Aquaculture* 13 (4):1958-1985.
393     Rhie, A., B.P. Walenz, S. Koren, and A.M. Phillippy, 2020 Merqury: reference-free quality,
394          completeness, and phasing assessment for genome assemblies. *Genome Biology* 21 (1):245.
395     Ripabelli, G., M.L. Sammarco, G.M. Grasso, I. Fanelli, A. Caprioli *et al.*, 1999 Occurrence of Vibrio and
396          other pathogenic bacteria in Mytilus galloprovincialis (mussels) harvested from Adriatic Sea,
397          Italy. *International Journal of Food Microbiology* 49 (1):43-48.
398     Roach, M.J., S.A. Schmidt, and A.R. Borneman, 2018 Purge Haplotigs: allelic contig reassignment for
399          third-gen diploid genome assemblies. *BMC Bioinformatics* 19 (1):460.
400     Rodríguez-Juíz, A.M., M. Torrado, and J. Méndez, 1996 Genome-size variation in bivalve molluscs
401          determined by flow cytometry. *Marine Biology* 126:489-497.
402     Ruan, J., and H. Li, 2020 Fast and accurate long-read assembly with wtdbg2. *Nature Methods* 17
403          (2):155-158.
404     Saco, A., B. Novoa, S. Greco, M. Gerdol, and A. Figueras, 2023 Bivalves Present the Largest and Most
405          Diversified Repertoire of Toll-Like Receptors in the Animal Kingdom, Suggesting Broad-
406          Spectrum Pathogen Recognition in Marine Waters. *Molecular Biology and Evolution* 40 (6).
407     Seuront, L., K.R. Nicastro, G.I. Zardi, and E. Goberville, 2019 Decreased thermal tolerance under
408          recurrent heat stress conditions explains summer mass mortality of the blue mussel Mytilus
409          edulis. *Scientific Reports* 9 (1):17498.
410     Smit, A., R. Hubley, and P. Green, 2015 RepeatMasker Open-4.0. *http://www.repeatmasker.org*.
411     Stanke, M., O. Keller, I. Gunduz, A. Hayes, S. Waack *et al.*, 2006 AUGUSTUS: ab initio prediction of
412          alternative transcripts. *Nucleic acids research* 34 (suppl_2):W435-W439.
413     Vaser, R., I. Sović, N. Nagarajan, and M. Šikić, 2017 Fast and accurate de novo genome assembly
414          from long uncorrected reads. *Genome Res* 27 (5):737-746.
415     Wang, Y., H. Tang, J.D. Debarry, X. Tan, J. Li *et al.*, 2012 MCScanX: a toolkit for detection and
416          evolutionary analysis of gene synteny and collinearity. *Nucleic acids research* 40 (7):e49.
417     Wilson, J., I. Matejusova, R.E. McIntosh, S. Carboni, and M. Bekaert, 2018 New diagnostic SNP
418          molecular markers for the Mytilus species complex. *PLoS One* 13 (7):e0200654.
419     Yang, J.L., D.D. Feng, J. Liu, J.K. Xu, K. Chen *et al.*, 2021 Chromosome-level genome assembly of the
420          hard-shelled mussel Mytilus coruscus, a widely distributed species from the temperate areas
421          of East Asia. *Gigascience* 10 (4).

422