1  **Genome assembly of the Australian black tiger shrimp (*Penaeus monodon*)**

2  **reveals a fragmented IHHNV EVE sequence**

3  Roger Huerlimann[*,†,‡,1,2], Jeff A Cowley[*,§,1], Nicholas M Wade[*,§,1], Yinan Wang[**], Naga

4  Kasinadhuni[**], Chon-Kit Kenneth Chan[**], Jafar Jabbari[**], Kirby Siemering[*,**], Lavinia

5  Gordon[**], Matthew Tinning[*,**], Juan D Montenegro[**,3], Gregory E Maes[††,‡‡], Melony J

6  Sellars[§,4], Greg J Coman[*,§§], Sean McWilliam[*,§], Kyall R Zenger[*,†], Mehar S Khatkar[*,***],

7  Herman W Raadsma[*,***], Dallas Donovan[*,†††], Gopala Krishna[*,†††], and Dean R Jerry[*,†,‡]

8  [*] ARC Industrial Transformation Research Hub for Advanced Prawn Breeding, Australia

9  [†] Centre for Sustainable Tropical Fisheries and Aquaculture, College of Science and

10  Engineering, James Cook University, Townsville, QLD 4811, Australia

11  [‡] Centre for Tropical Bioinformatics and Molecular Biology, James Cook University,

12  Townsville, QLD 4811, Australia

13  [§] CSIRO Agriculture and Food, 306 Carmody Road, St Lucia, QLD, 4067, Australia

14  [**] Australian Genome Research Facility Ltd, Level 13, Victorian Comprehensive Cancer

15  Centre, 305 Grattan St, Melbourne VIC 3000, Australia

16  [††] Laboratory of Biodiversity and Evolutionary Genomics, KU Leuven, Leuven, 3000,

17  Belgium

18  [‡‡] Center for Human Genetics, UZ Leuven- Genomics Core, KU Leuven, Leuven, 3000,

19  Belgium

20  [§§] CSIRO Agriculture and Food, Integrated Sustainable Aquaculture Program, 144 North

21  Street, Woorim, QLD 4507, Australia

22  [***] Sydney School of Veterinary Science, Faculty of Science, The University of Sydney,

23  Camden, NSW 2570, Australia

24  [†††] Seafarms Group Ltd, Level 11 225 St Georges Terrace, Perth, WA 6000, Australia

25  [1] Shared first authors

26  [2] Corresponding author: Roger Huerlimann; James Cook University, 145 James Cook

27  Street, Townsville, QLD 4811, Australia; roger.huerlimann@jcu.edu.au

28    [3] Current address: Genics Pty Ltd, Level 5, 60 Research Road, St Lucia, QLD, 4067,

29    Australia

30    [4] Current address: Department of Neurosciences and Developmental Biology,

31    University of Vienna, Vienna BioCenter, Vienna 1030, Austria

32    *Keywords:* *Penaeus monodon*, Australia, genome assembly, PacBio, IHHNV EVE

33    Running title: **Australian black tiger shrimp genome**

34    **<u>Abstract</u>**

35    Shrimp are a valuable aquaculture species globally; however, disease remains a major

36    hindrance to shrimp aquaculture sustainability and growth. Mechanisms mediated by

37    endogenous viral elements (EVEs) have been proposed as a means by which shrimp

38    that encounter a new virus start to accommodate rather than succumb to infection over

39    time. However, evidence on the nature of such EVEs and how they mediate viral

40    accommodation is limited. More extensive genomic data on Penaeid shrimp from

41    different geographical locations should assist in exposing the diversity of EVEs. In this

42    context, reported here is a PacBio Sequel-based draft genome assembly of an

43    Australian black tiger shrimp (*Penaeus monodon*) inbred for one generation. The 1.89

44    Gbp draft genome is comprised of 31,922 scaffolds (N50: 496,398 bp) covering 85.9%

45    of the projected genome size. The genome repeat content (61.8% with 30%

46    representing simple sequence repeats) is almost the highest identified for any species.

47    The functional annotation identified 35,517 gene models, of which 25,809 were protein-

48    coding and 17,158 were annotated using interproscan. Scaffold scanning for specific

49    EVEs identified an element comprised of a 9,045 bp stretch of repeated, inverted and

50    jumbled genome fragments of Infectious hypodermal and hematopoietic necrosis virus

51    (IHHNV) bounded by a repeated 591/590 bp host sequence. As only near complete

52    linear ~4 kb IHHNV genomes have been found integrated in the genome of *P. monodon*

53    previously, its discovery has implications regarding the validity of PCR tests designed to

54    specifically detect such linear EVE types. The existence of joined inverted IHHNV

55    genome fragments also provides a means by which hairpin dsRNAs could be expressed

56    and processed by the shrimp RNA interference (RNAi) machinery.

57                                          **<u>INTRODUCTION</u>**

58    Shrimp aquaculture plays a central role in producing high quality protein for human

59    consumption, with global aquaculture production of the two major species, *Penaeus*

60    *vannamei* and *P. monodon*, reaching close to six million tons in 2018 (FAO 2020).

61    However, diseases, such as those caused by highly pathogenic viruses, are currently a

62    major contributor to unfulfilled production potential (FAO 2020). Therefore, a more

63    advanced understanding of the host defense mechanisms that suppress infection will be

64    critical to finding solutions to viral diseases (Kulkarni et al. 2021; Hauton 2017; Yang et

65    al. 2021).

66    Initially described in insects, the viral accommodation mechanism has been

67    hypothesized to explain why farmed shrimp highly susceptible to morbidity and mortality

68    proceeding their initial encounter with a new virus tend to become less susceptible over

69    time (Flegel 2020). Viral accommodation is mediated through host-genome integrated

70    endogenous viral elements (EVEs) that can be inherited after integration into the germ

71    line. The expressed EVE-specific double-stranded RNA (dsRNA) is then processed by

72    the host RNA interference (RNAi) pathway, suppressing viral RNA expression levels

73    and therefore infection loads. In the case of RNA viruses, a linear copy viral DNA

74    (cvDNA) or circular copy viral DNA (ccvDNA) can be reverse transcribed by the host

75    (Taengchaiyaphum et al. 2021). These DNA copies of virus RNA can then either

76    autonomously insert into the host genome to become an EVE, or be used directly as a

77    template for dsRNA transcription as an initial step to RNAi-mediated suppression of

78    virus infection (Taengchaiyaphum et al. 2021).

79    Of the >50,000 known crustacean species, high-quality genome assemblies are only

80    available for a select few taxa, driven primarily by the commercial or unique biological

81    significance of certain species. Genome assemblies provide a reference base for

82    functional transcriptomic studies (Yue and Wang 2017; Chandhini and Rejish Kumar

83    2019), aid in the positioning of genetic markers used for selective breeding (Houston et

84    al. 2020; Zenger et al. 2017) and provide an important resource for the examination and

85    characterization of genomic regions of commercial or biological interest (Guppy et al.

86    2020; Hollenbeck and Johnston 2018). However, crustacean genomes have also

87  proved immensely challenging to assemble due to their large (>2 Gbp), highly repetitive

88  (>50%), and highly heterozygous genomes (Yuan et al. 2021a). To some extent, these

89  difficulties have been alleviated by the advent of single-molecule long-read sequencing

90  and improved genome assemblers. Extracting intact high-quality genomic DNA from

91  muscle tissue of crustaceans like shrimp has also proved problematic and exacerbated

92  difficulties in obtaining high-quality data from various NGS platforms (Angthong et al.

93  2020). Despite these challenges, genome assemblies highly fragmented into more than

94  a million contigs have been reported for the penaeid shrimp species *P. vannamei* (Yu et

95  al. 2015), *P. japonicus* (Yuan et al. 2018), and *P. monodon* (Van Quyen et al. 2020;

96  Yuan et al. 2018). Through applying long-read sequencing and HiC scaffolding, less

97  fragmented high-quality genomes have also been achieved recently for *P. vannamei*

98  (Zhang et al. 2019), *P. monodon* (pseudo-chromosome level) (Uengwetwanit et al.

99  2020) and *P. japonicus* (Kawato et al. 2021).

100  Reported here is a high-quality draft genome assembly of a single-generation inbred

101  male *P. monodon* from eastern Australia, a population genetically distinct from others

102  across their South East Asian, Indo-Pacific and East African distribution (Vu et al.

103  2021). We report and resolve the genomic structure of an EVE of Infectious hypodermal

104  and hematopoietic necrosis virus (IHHNV) comprised of repeated, inverted, and jumbled

105  IHHNV genome fragments. We discuss the disease detection implications of false PCR-

106  positives for infectious IHHNV, and how the EVE might have originated.

## METHODS & MATERIALS

**Shrimp breeding and selection for sequencing:**

109  A second-generation (G2) male *Penaeus monodon* that had undergone a single cycle of

110  inbreeding was selected for genomic sequencing. The original wild-caught broodstock

111  were collected from a Queensland east coast location (approximately 17.3°S, 146.0°E)

112  in September 2013. In October 2013, 14 first-generation (G1) families were produced

113  from the brood stock at Seafarm Flying Fish Point hatchery (approximately 17.5°S,

114  146.1°E). In February 2015, pleopod tissue was sampled from 50 female and 50 male

115  G1 broodstock. These tissues were genotyped (using 2 x 60 SNP panels (Sellars et al.

116  2014) to identify the parental origin of each broodstock and to select related mating

117   pairs to generate the inbred G2 progeny. In August 2015, groups of 50 juvenile males

118   from 5 inbred G2 families were euthanized to collect muscle tissue from the first

119   abdominal segment for sequencing and the second most anterior pair of pleopods for

120   genotyping. These tissues, as well as the remainder of each shrimp (archived source of

121   tissue for sequencing) were snap frozen under dry ice pellets and stored at -80°C. Each

122   shrimp was then genotyping using the 120-SNP panel (Sellars et al. 2014) and a

123   genome-wide SNP assay based on DArTSeq (Guppy et al. 2020). After ranking the 50

124   males based on inbreeding coefficient (F) and multi-locus heterozygosity (MLH) data

125   from the 120-SNP panel, the individual (named Nigel) with the highest inbreeding

126   coefficient was chosen for genomic sequencing. The choice was confirmed using a

127   genome-wide SNP assay based on DArTSeq of the top five inbred shrimp based on the

128   120-SNP panel which recovered the same ranking (Nigel: MLH of 0.231 and F of

129   0.271).

130   **DNA extraction, library preparation and genome sequencing:**

131   Multiple extraction methods were trialed to generate intact high-quality genomic DNA

132   from stored muscle tissue of the single selected inbred shrimp. All DNA extractions and

133   sequencing runs were carried out at the Australian Genome Research Facility (AGRF),

134   Melbourne, Australia. For Illumina sequencing, the MagAttract HMW DNA kit (QIAGEN)

135   was used and PCR-free fragment shotgun libraries were prepared using the 'with-bead

136   pond library' construction protocol described by Fisher et al. (Fisher et al. 2011) with

137   some modifications (Supplementary Material 1). The library was sequenced on two

138   HiSeq 2500 lanes using a 250 bp PE Rapid sequencing kit (Illumina). The same DNA

139   was also used to create a 10X Genomics Chromium library as per manufacturer

140   instructions, which was sequenced on two HiSeq 2500 lanes using a 250 bp PE Rapid

141   sequencing kit. For PacBio sequencing, the following DNA extraction methods were

142   used with varying success: MagAttract HMW DNA kit (QIAGEN), Nanobind HMW

143   Tissue DNA kit-alpha (Circulomics), and CTAB/Phenol/Chloroform (Supplementary

144   Table 1). Libraries were prepared using the SMRTbell Template Prep Kit 1.0 (PacBio),

145   loaded using either magbeads or diffusion, and sequenced using the Sequel

146   Sequencing Kits versions 2.1 and 3.0 on a PacBio Sequel (Supplementary Table 1).

147   The same muscle tissue was also used to prepare three Dovetail Hi-C libraries

148    according to manufacturer's instructions. Two libraries were sequenced on a shared

149    lane of a NovaSeq S1 flow cell, and a third library was sequenced on one lane of a

150    NovaSeq SP flow cell, with both sequencing runs generating 100 bp paired-end reads.

**Genome assembly:**

152    The quality of the initial short-read genome assemblies using either DISCOVAR *de*

153    *novo* (Weisenfeld et al. 2014) with Illumina data, or Supernova (Weisenfeld et al. 2017)

154    with 10X Genomics Chromium data was poor. The most contiguous assembly was

155    achieved using wtdbg2/redbean (Version 2.4, Ruan and Li 2019) with 75 X times

156    coverage of PacBio data, setting the estimated genome size to 2.2 Gb, but without

157    using the wtdbg2 inbuilt polishing. The raw assembly was subjected to two rounds of

158    polishing using the PacBio subreads data in arrow (Version 2.3.3,

159    github.com/PacificBiosciences/GenomicConsensus) and one round of polishing using

160    the Illumina short-read data in pilon (Version 1.23, Walker et al. 2014). Scaffolds were

161    constructed in two steps. Medium-range scaffolding carried out using 10X Genomics

162    Chromium data with longranger (Version 2.2.2,

163    https://support.10xgenomics.com/genome-exome/software/downloads/latest) and

164    ARCS (Version 1.0.6, Yeo et al. 2017), while long-range scaffolding was performed

165    using dovetail Hi-C data, and intra- and inter-chromosomal contact maps were built

166    using HiC-Pro (Version 2.11.1, Servant et al. 2015) and SALSA (commit version

167    974589f, Ghurye et al. 2017). This genome assembly was then submitted to NCBI

168    GenBank, which required the removal of two small scaffolds and the splitting of one

169    scaffold. The overall quality of the final V1.0 genome was assessed using BUSCO, and

170    through mapping of RNA-seq, and Illumina short-reads using HiSAT2 (version 2.1.0,

171    Kim et al. 2019).

**Repeat annotation:**

173    Repeat content was assessed with *de novo* searches using RepeatModeler (V2.0.1)

174    and RepeatMasker (V4.1.0) via Dfam TE-Tools (V1.1, https://github.com/Dfam-

175    consortium/TETools) within Singularity (V2.5.2, Kurtzer et al. 2017). Additionally,

176    tandem repeat content was determined using Tandem Repeat Finder (V4.0.9, Benson

177    1999) within RepeatModeler. Analyses and plotting of interspersed repeats were carried

178    out as per Cooke *et al.* (2020,

179    github.com/iracooke/atenuis_wgs_pub/blob/master/09_repeats.md). Additionally, the

180    genomes of the Black tiger shrimp (Thai origin, www.biotec.or.th/pmonodon; Kim et al.

181    2019), Whiteleg shrimp (*P. vannamei*, NCBI accession: QCYY00000000.1; Zhang et al.

182    2019), Japanese blue crab (*Portunus trituberculatus*, gigadb.org/dataset/100678; Tang

183    et al. 2020), and Chinese mitten crab (*Eriocheir japonica sinensis*, NCBI accession:

184    LQIF00000000.1) were run through the same analyses for comparison.

185    **Gene prediction and annotation:**

186    In order to generate an RNA-seq based transcriptome, raw data from a previous study

187    (NCBI project PRJNA421400; Huerlimann et al. 2018) was mapped to the masked

188    genome using STAR (Version 2.7.2b; Dobin et al. 2013), followed by Stringtie (Version

189    2.0.6; Pertea et al. 2015) (Supplementary Table 2). Additionally, the IsoSeq2 pipeline

190    (PacBio) was used to process the ISO-seq data generated in this study (Supplementary

191    Table 2). Finally, the genome annotation was carried out in MAKER2 (v2.31.10;

192    Campbell et al. 2014; Cantarel et al. 2008; Holt and Yandell 2011) using the assembled

193    RNA-seq and ISO-seq transcriptomes together with protein sequences of other

194    arthropod species (Supplementary Table 3).

195    **Endogenous viral element analysis:**

196    BLASTn using a 3,832 bp IHHNV EVE Type A sequence detected in Australian *P.*

197    *monodon* (Au2005; EU675312.1) as a query identified a potential EVE in Scaffold_97 of

198    the *P. monodon* genome assembly. The EVE was unusual in that it comprised of

199    repeated, inverted and jumbled fragments of an EVE Type A sequence. The nature and

200    arrangement of EVE fragments was initially determined manually and the relative

201    sequence positions of matching fragments within the EVE and scaffold sequence was

202    determined using QIAGEN CLC Genomics Workbench 18.0

203    (https://digitalinsights.qiagen.com/). To confirm the authenticity of the Scaffold_97 EVE

204    (S97-EVE), six PCR primer sets were designed using Primer 3 v.0.4.0 (Koressaar and

205    Remm 2007; Untergasser et al. 2012) to amplify each EVE boundary and two internal

206    sequences (Supplementary Table 4). DNA was extracted from ~10 mg gill tissue stored

207    at -80°C from the *P. monodon* sequenced using DNAeasy kit spin columns (QIAGEN).

208 DNA was eluted in 50 µL EB buffer, aliquots were checked to DNA concentration and

209 purity using a Nanodrop 8000 UV spectrophotometer and the remainder was stored at -

210 20°C. As DNA yields were low (9-38 ng/µL), a 1.0 µL aliquot of each sample was

211 amplified in 10 µL reactions incubated at 30°C for 16 h as described in the REPLI-g Mini

212 Kit (QIAGEN).

213 Each PCR (25 µL) contained 2 µL REPLI-g amplified gill DNA, 1 x MyTaq$^{TM}$ Red Mix

214 (Bioline), 10 pmoles each primer and 0.25 µL (1.25 U) MyTaq DNA Polymerase

215 (Bioline). Thermal cycling conditions were 95°C for 1 min followed by a 5-cycle touch-

216 down (95°C for 30 s, 60°C to 56°C for 30 s, 72°C for 20 s), 30 cycles of the same using

217 an anneal of 55°C for 30 s, followed by 72°C for 7 min and a 20°C hold. For semi-

218 nested PCR using the 1b and 4b primer sets, 1 µL each PCR (either neat or diluted 1:5

219 to 1:10 depending on PCR product amount) was amplified similarly for 30 cycles using

220 an anneal step of 55°C for 30 s. Aliquots (5-10 µL) of each reaction were

221 electrophoresed in a 1.0% agarose-TAE gel containing 0.1 µL mL$^{-1}$ ethidium bromide,

222 and a gel image was captured using a Gel Doc 2000 UV transilluminator (Bio-Rad).

223 Each amplicon was purified using a spin column (QIAGEN) and sequenced at the

224 Australian Genome Research Facility (AGRF), Brisbane. The quality of sequence

225 chromatograms was evaluated and consensus sequences for each amplicon were

226 generated using Sequencher® 4.9 (Gene Codes Corp.).

227 **Data availability:**

228 Raw and assembled sequence data generated by this study have been deposited in

229 GenBank BioProject PRJNA590309, BioSample SAMN13324362. PacBio and Illumina

230 raw data can be found under accession numbers SRR10713990-SRR10714025. The

231 final scaffolded assembly can be found under accession JAAFYK000000000. RNA-seq

232 data used for annotation originated from an earlier study (Huerlimann et al. 2018).

233 **RESULTS AND DISCUSSION**

234 **DNA extraction, library preparation and genome sequencing:**

235 In total, 158 Gb (72 X coverage) of Illumina, 494 Gb (224 X coverage) of 10X Genomics

236 Chromium, 165 Gb (75 X coverage) of PacBio Sequel, and 119 Gb (54 X coverage) of

237     DoveTail data were generated (Table 1). While the MagAttract HMW DNA kit (QIAGEN)

238     was suitable for Illumina sequencing (PCR-free shotgun libraries and 10X Genomics

239     Chromium), using this DNA resulted in poor PacBio Sequel sequencing runs

240     (Supplementary Table 1). Runs consistently showed low yield and short fragment

241     lengths, despite relatively high molecular weight DNA. However, DNA extracted with the

242     Nanobind HMW Tissue DNA kit-alpha (Circulomics, Inc., Baltimore, USA) showed better

243     sequencing performance (higher yield and fragment length; Supplementary Table 1).

244     Furthermore, diffusion loading of the PB Sequel resulted in better results than magbead

245     loading. DNA derived from either extraction method was unsuitable for Oxford

246     Nanopore Technology (ONT) sequencing due to it rapidly blocking the pores (data not

247     shown).

248     Sequence quality issues associated with DNA extraction have also been noted in other

249     shrimp genome assembly reports (Zhang et al. 2019; Uengwetwanit et al. 2020). The

250     patterns seen in the PacBio sequencing results (short polymerase read lengths despite

251     high quality libraries), coupled with the inability to successfully sequence *P. monod*on

252     using ONT technology (immediate pore blockage), can be explained by high amounts of

253     polysaccharides and polyphenolic proteins co-extracting with the DNA. This has also

254     been mentioned by Angthong et al. (2020), who also present an alternative DNA

255     extraction method to the Circulomics Nanobind HWM Tissue DNA extraction kit

256     suggested here.

**Genome assembly and quality assessment:**

258     As reported by other Penaeid shrimp genome sequencing projects (Uengwetwanit et al.

259     2020; Zhang et al. 2019; Yuan et al. 2021a), sequencing and assembly of the Australian

260     *P. monodon* genome proved problematic due to its large size, substantial

261     heterozygosity and prevalence of repeat elements. The *de novo* assembly of the PacBio

262     data resulted in 47,607 contigs (contig N50: 77,900 bp) a total of 1.90 Gbp in size

263     (Table 2). After medium-range scaffolding with 10X Genomic Chromium data and long-

264     range scaffolding with Dovetail sequences, the resulting scaffolded assembly contained

265     1.89 Gbp across 31,922 scaffolds (scaffold N50: *496*,398 bp; Table 2). Assuming a

266     genome size of 2.2 Gbp (Huang et al. 2011), this scaffolded assembly covers 85.9 % of

267  the projected *P. monodon* genome (Table 2). This is slightly lower than the 90.3%

268  recently achieved for the same species in Thailand (Uengwetwanit et al. 2020), and

269  higher than the 67.7% achieved for *P. vannamei* (Zhang et al. 2019), which has a

270  slightly larger genome. Altogether, 98.1% of the Illumina DNA short-read data mapped

271  to the raw assembly. BUSCO (V3; Simão et al. 2015), using the Arthropoda odb9

272  database (Zdobnov et al. 2017), estimated the Australian *P. monodon* genome

273  assembly to be 86.8% complete (gene n = 1,066; 85.8% single copy; 1.0% duplicated;

274  4.5% fragmented; 8.7% missing; Table 2). These assembly metrics are comparable to

275  those achieved for the Thai *P. monodon* assembly (C 87.9%, S 84.8%, D 3.1%, F 4.0%,

276  M 8.0%; Uengwetwanit et al. 2020) and slightly better than those achieved for the *P.*

277  *vannamei* assembly (C 78.0%, S 74.0%, D 4.0%, F 4.0%, M 18.0%; Zhang et al. 2019),

278  both analyzed with the same database and BUSCO version (Table 2).

279  **Functional and repeat annotation:**

280  The functional annotation using RNA-seq, ISO-seq and protein information, identified

281  35,517 gene models, of which 25,809 were protein-coding and 17,158 were annotated

282  using interproscan (Table 2). Similar numbers of genes were annotated in the Thai *P.*

283  *monodon* (Uengwetwanit et al. 2020) and *P. vannamei* (Zhang et al. 2019) assemblies.

284  Repeat content in the Australian *P. monodon* assembly (61.8%) was high, like in the

285  Thai *P. monodon* assembly (62.5%; Uengwetwanit et al. 2020), and substantially higher

286  than in genome assemblies of *P. vannamei* (51.7%; Zhang et al. 2019)), *Portunus*

287  *trituberculatus* (45.9%, Tang et al. 2020) or *Eriocheir japonica sinensis* (35.5%,

288  LQIF00000000.1) (Supplementary Table 5, Fig. 1). Interestingly, simple sequence

289  repeats (SSRs) that dominated in prevalence (30.0%) in the Australian *P. monodon*

290  assembly were less prevalent (23.9%) in the Thai *P. monodon* assembly (Uengwetwanit

291  et al. 2020), similarly prevalent (27.1%) in the *P. vannamei* assembly, but far less

292  prevalent in the genome assemblies of either the Japanese blue (16.9%) or Chinese

293  mitten crab (7.9%) (Supplementary Table 5, Fig. 1). Such high SSR levels have been

294  linked to genome plasticity and adaptive evolution facilitated through transposable

295  elements (Yuan et al. 2021b). In addition to SSRs, the Australian *P. monodon* assembly

296  contained 9.8% long interspersed nuclear elements (LINEs), 2.5% low complexity

297  repeats, 2.0% DNA transposons, 1.6% long terminal repeats (LTRs), 0.51% small

298   interspersed nuclear elements (SINEs), 0.1% satellites, 0.01% small RNA repeats and

299   15.4% unclassified repeat element types (Supplementary Table 5, Fig. 1). Broad

300   comparisons of the major repeat types in the genome assemblies of *P. monodon*, *P.*

301   *vannamei*, *Portunus trituberculatus* and *E. japonica sinensis* based on kimura distances

302   showed them to be relatively conserved across all four crustacean species (Fig. 1). At

303   all lengths and levels of divergence, unknown repeats dominated, with a large

304   proportion of these >100 kbp in size (Fig. 1A). Repeat patterns shared across the four

305   species were further highlighted when unknown reads were removed, and repeats split

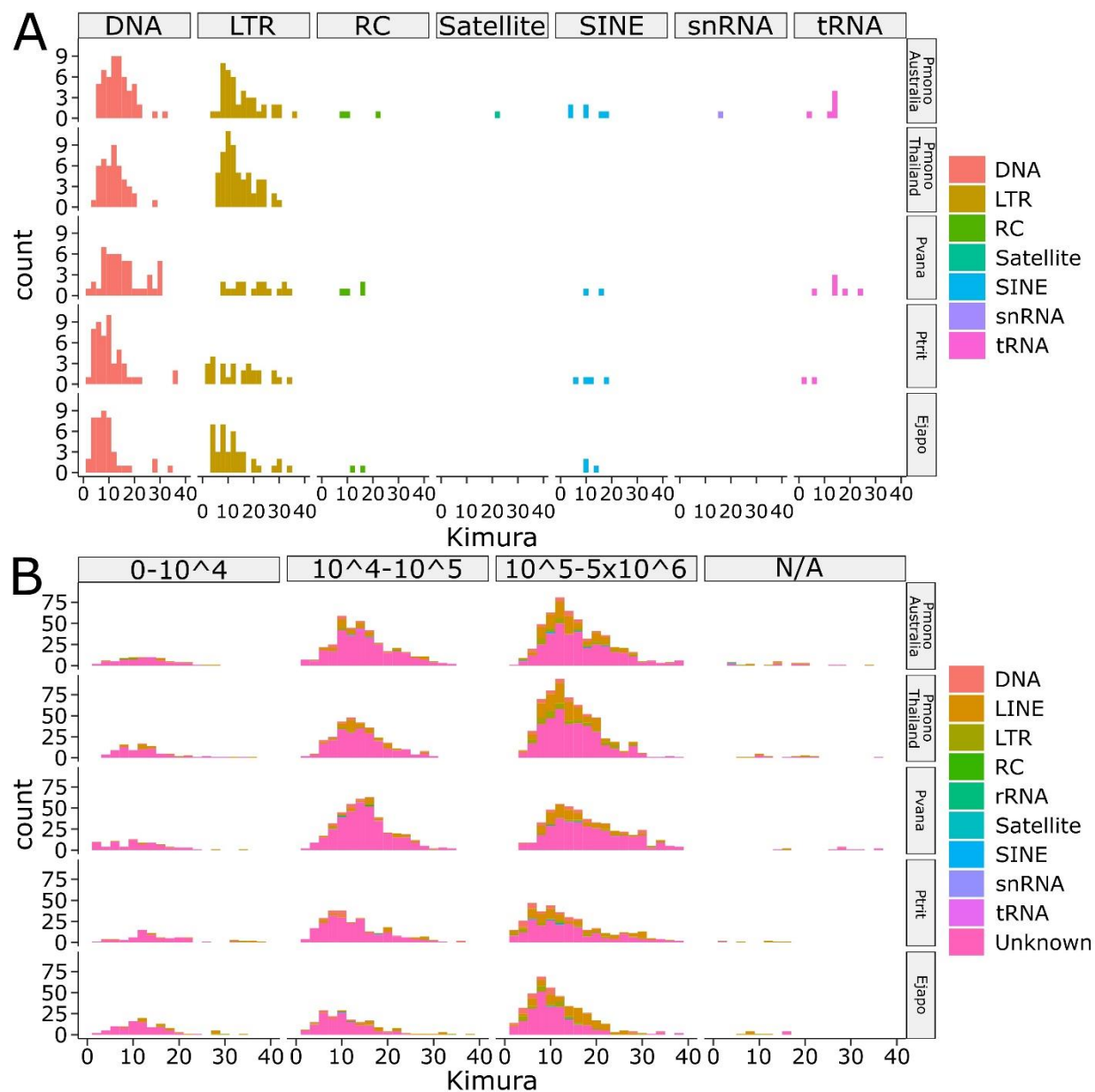306   into major classes (Fig. 1B).

**Fig. 1** Kimura distances of repetitive sequences in the genome assemblies of Australian black tiger shrimp (*P. monodon*, NCBI accession: JAAFYK000000000, Pmono Australia, this study) Thai black tiger shrimp (Pmono Thailand, *P. monodon*, Pmono Thailand, Uengwetwanit et al. 2020), Whiteleg shrimp (Pvana, *Penaeus vannamei*, NCBI accession: QCYY00000000.1, Zhang et al. 2019), Japanese blue crab (Ptrit, *Portunus trituberculatus*, gigadb.org/dataset/100678, Tang et al. 2020), and Chinese mitten crab (Ejapo*, Eriocheir japonica sinensis*, NCBI accession: LQIF00000000.1) determined by using either (A) repeat length or (B) repeat class.

**IHHNV-EVE rearrangement in the Australian *P. monodon* genome:**

Sequences homologous to a 3,832 bp linear IHHNV-EVE (Au2005, Type A) found to occur in some Australian *P. monodon* (Krabsetsve et al. 2004) were identified in Scaffold_97 (S97, 2,608,951 nt). However, rather than representing an intact linear copy of this EVE, the S97-EVE comprised a 9,045 bp stretch of jumbled, repeated, and inverted IHHNV fragments flanked by two repeated 591/590 bp (flanking repeat) sequences (Fig. 2). Alignments identified most fragments to be jumbled relative to their location in the Au2005 IHHNV-EVE sequence, and the expanded EVE length to be due to replicated short sequences originating from 5'-terminal genome regions. Fragments positioned at the S97-EVE extremities generally originated from the central and downstream regions of the Au2005 IHHNV-EVE sequence and were consistently orientated inwards. The central S97-EVE region comprised a block of at least six 661 bp repeat units (RUs). Each RU was comprised of two inward-facing sequences either (A) 398 bp or (B) 263 bp in length that mapped to the same region (94-501 and 94-368, respectively) at the 5'-terminus of the Au2005 IHHNV-EVE (Fig. 2B, grey arrows). In total, 83% of the Au2005 IHHNV-EVE sequence was identified to be covered by genome fragments present in the S97-EVE, with those present being on average 99.3% identical.

The inverted A and B sequences comprising each RU contain RNA transcription regulatory signals of the IHHNV P2 promoter (Shike et al. 2000; Dhar et al. 2011; Dhar et al. 2010; Dhar et al. 2007). Both initiated at a sequence (5'-GTCATAGGT…) mapping precisely to a G nucleotide residing immediately downstream of the inversion point (|) of an 18 bp inverted repeat (5'-..TTACAACCTATGAC|GTCATAGGTCCTATATAAGAGT..-3') located 2 bp upstream of the TATA-box element (5'-TATATAA-3') of the P2 transcriptional promoter (Dhar et al. 2011; Dhar et al. 2010; Dhar et al. 2007). The A and B repeat components in each RU of the six blocks were orientated 5'|B-A|B-A|B-A|A-B|B-A|B-A|3', with those in RU4 being reversed compared to the others. Due to the A and B repeat components being inverted, the 18 bp inverted repeat (ie. 5'-..ACTCTTATATAGGACCTATGAC|GTCATAGGTCCTATATAAGAGT..-3') was reconstructed at each of the 5 RU junction sites irrespective of which 2 repeat components (A|A, A|B or B|B) were joined (Fig. 2B, purple bars). This arrangement

347     generated a 544 bp inverted repeat (263 x 2 + 18) for sequences extending from either

348     A|B or B|B RU junctions, or a 1,902 bp inverted repeat (661 x 2 + 263 x 2 + 18 x 3) for

349     the long complimentary sequence stretches extending outwards from the A|A

350     components at the RU3|RU4 junction to the end of repeat component A of RU2 and the

351     equivalent position of repeat component B in RU5. However, relating to the descriptions

352     of this unusual EVE segment, it is important to note that no single long read was

353     obtained that traversed the entire six RU blocks into flanking unique S97-EVE

354     sequences (Fig. 2). Combined with short read numbers generated using various

355     sequencing methods being substantially elevated at positions mapping to each block

356     RU (Fig. 2C), the likelihood of the block comprising more than six RUs remains to be

357     established.

358     DNAFold and RNAfold analyses showed the 18 bp inverted repeat, the inverted A and B

359     repeat components of each RU and the longer complimentary sequences that stretched

360     through multiple RUs to all have potential to form highly stable simple to complex

361     secondary structures as either ssDNA or ssRNA (data not shown). Discrete DNA

362     secondary structures are known to have roles in mediating recombination in mobile

363     genetic elements (Bikard et al. 2010) and in the genomes of parvoviruses like the

364     extensively studied adeno-associated virus (AAV), structures formed by inverted

365     terminal repeat (ITR) sequences play critical roles in initiating genomic ssDNA

366     replication, genomes forming circular extrachromosomal dsDNA episomes and genomic

367     integrating into host chromosomal DNA (Cotmore and Tattersall 1996, Kotin et al. 1991,

368     Schnepp et al. 2005, Yang et al. 1997). The mechanisms leading to the A and B

369     inverted repeat sequences forming the 661 bp RUs and their apparent multiplication in

370     the central region of the S97-EVE remains unknown. However, their existence is

371     consistent with integrated AAV proviral DNA structures being observed to contain head-

372     to-tail tandem arrays of partial ITR sequences and for genomic rearrangements

373     occurring via deletion and/or rearrangement-translocation at the integration site (Yang

374     et al. 1997).

375     The 18 bp inverted repeat at the S97-EVE RU junctions also occurred at the upstream

376     RU1 and downstream RU6 boundaries of the RU block. However, unlike those at the

377     internal RU junctions which extended into the same downstream Au2005-EVE

378    sequence including the TATA-box element (Dhar et al. 2011; Dhar et al. 2010; Dhar et

379    al. 2007; Krabsetsve et al. 2004), the outer half of each inverted repeat flanking the RU-

380    block extended into sequences toward the 5' end of the IHHNV genome

381    (Supplementary Figure 1). Three disparate partial RU sequences (pRUa, pRUb, pRUc)

382    associated with four 18 bp inverted repeats also resided just upstream of the 6 RU

383    block. Like RU1 and RU6, one side of each inverted repeat possessed variable lengths

384    of sequence extending toward the IHHNV genome 5'-terminus (Supplementary Figure

385    1).

386    In some IHHNV strains, the sequence immediately upstream of the 18 bp inverted

387    repeat comprises a second imperfect 39-40 bp inverted repeat. With an IHHNV strain

388    detected in Pacific blue shrimp (*Penaeus stylirostris*) sampled from the Gulf of California

389    in 1998 (Shike et al. 2000, AF273215.1), the 5'-genome terminus upstream of it

390    consisted of an 8 bp portion of the downstream 18 bp inverted repeat (Supplementary

391    Figure 1). In the S97-EVE, the 18 bp inverted repeats associated with each terminal RU

392    or upstream pRU extended 18-38 bp into the 39-40 bp inverted repeat (Supplementary

393    Figure 1). Of interest, with the first pRU occurring in the S97-EVE (5'-pRUa), the 93 bp

394    sequence abutting the 18 bp inverted repeat was also identical to the 5'-terminal

395    sequence reported for the Au2005 IHHNV-EVE found in *P. monodon* sampled from

396    farms in Australia in 1993/1997 (Krabsetsve et al. 2004; EU675312.1).

397    To confirm that the fragmented and jumbled nature of the S97-EVE was not an

398    assembly artefact, regions spanning each EVE extremity to unique host sequences

399    positioned just beyond the 591/590 bp flanking repeats, as well as two internal regions

400    each spanning conjoined non-repeated EVE fragments were amplified by PCR

401    (Supplementary Table 4; Fig. 2D). Amplicons of the expected sizes were clearly

402    amplified by each extremity PCR test (S97-1a and S97-4a) and the S97-3 internal PCR

403    test (Fig. 2D). The other internal PCR test (S97-2) also generated a 1,337 bp amplicon

404    of the expected size, as well as one ~200 bp shorter, but in relatively lower abundance.

405    Using each extremity PCR product as template, semi-nested PCR tests using an

406    alternative internal EVE-specific primer also produced amplicons of the expected

407    shorter sizes, and their authenticity was confirmed by sequence analysis (data not
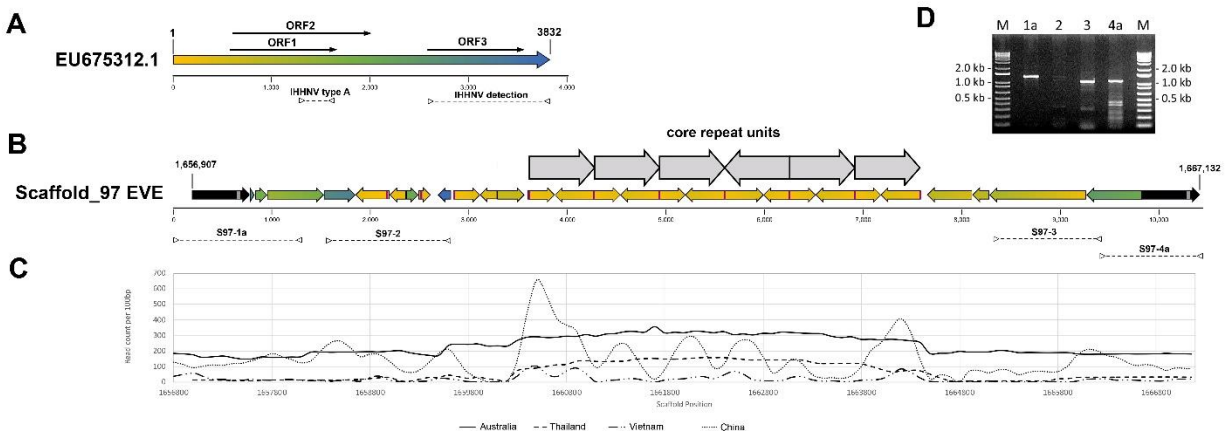
408    shown).

409

**Fig. 2** (A) Schematic diagram of a 3,832 bp ssDNA genome of Infectious hypodermal and hematopoietic necrosis virus (IHHNV) showing the relative positions of coding sequences (arrows) for the virus replicase (ORF1), NS1 non-structural protein (ORF2) and viral capsid protein (ORF3). A colour gradient was applied to visualize relative genome positions. (B) Schematic diagram of the positions and orientations of IHHNV genome fragments comprising the Scaffold_97 EVE (S97-EVE). The orientations of the IHHNV fragments (coloured arrows) and the flanking repeated 591/590 bp host sequence (black arrows) are shown by arrow directions. The origins of the S97-EVE fragments relative to their positions in a linear IHHNV-EVE (see A) are identified by colour. The 10,226 bp S97-EVE resided between positions 1,656,907 and 1,667,132 in the 2,608,951 bp Scaffold_97 sequence. The larger grey arrows identify the positions and orientations of at least 6 core repeat blocks comprising of 2 smaller inverted repeats. Grey vertical bars show the location of a 34 bp sequence in each flanking repeat capable of folding into a stable secondary structure. The purple vertical bars show the locations of the 18 bp palindromic sequence present at the boundaries of each repeat unit (RU) and partial RU. Dashed lines (>--<) identify the regions amplified by the 4 PCR tests S97-1a, S97-2, S97-3, and S97-4a. (C) Coverage depth across the S97-EVE sequence of raw short reads used to assemble genome scaffolds of *P. monodon* from Australia (this study), Thailand (Uengwetwanit et al. 2020), Vietnam (Van Quyen et al. 2020) and China (Yuan et al. 2018). (D) Agarose gel image showing DNA products amplified by the S97-1a, S97-2, S97-3, and S97-4a PCR tests.

*P. monodon* **repeat sequences flanking the IHHNV-EVE:**

432  BLASTn and BLASTx searches did not identify any homologues of the 591/590 bp

433  flanking repeat sequence in GenBank. However, searches of the *P. monodon* genome

434  assembly identified long closely-related sequences in hundreds of other scaffolds (data

435  not shown). The searches also highlighted the presence of a 34 bp sequence

436  (5'-..ATGACTCCTCCCCCATAGATAGGGGCGGAGTCAT..-3') in each flanking repeat

437  (Fig. 2B, grey bars, upstream repeat position 1,657,364-1,657,397; downstream repeat

438  position 1,667,000-1,667,033) that was also present in 178 other scaffolds at >80%

439  identity. DNAFold and RNAfold analyses showed the sequence and its reverse

440  compliment to fold into stable hairpin structures as either ssDNA (ΔG = -10.44/-11.92,

441  Tm = 83.8/85.7℃) or ssRNA (ΔG = -20.40/-23.70). However, whether this or other

442  sequences in the host flanking repeat interact with IHHNV genome sequences and

443  proteins to facilitate recombination and site-specific integration remains to be

444  investigated. In this regard, the flanking host repeat possessed a

445  5'..CTTACTTACACTTG..3' tetramer repeat, which to the 5'-side of the S97-EVE was

446  located 33 bp upstream of the IHHNV CTTA.. sequence at the host/S97-EVE junction,

447  much like the host tetramer repeats well characterised to be pivotal to the AAV genome

448  integrating at a specific location in human chromosome 19 (Kotin et al. 1992, Linden et

449  al. 1996).

450  **Comparison to jumbled IHHNV-EVEs in other *P. monodon* genome assemblies:**

451  BLASTn searches of the most comprehensive genome assembly of a *P. monodon* from

452  Thailand (NSTDA_Pmon_1, GCA_015228065.1, Uengwetwanit et al. 2020) identified

453  Scaffold_35 (S35) containing two disparate aggregations of jumbled IHHNV-EVE Type

454  A fragments (S35-EVE1 = 7,888 bp; S35-EVE2 = 16,310 bp) each flanked by >500 bp

455  host repeats near identical in sequence to those flanking the S97-EVE (Table 3).

456  Compared to the S97-EVE, 2,328 bp of S35-EVE1 sequence immediately downstream

457  of the 5' 592 bp host repeat, except for a 166 bp deletion, and 647 bp of sequence

458  immediately upstream of the 3' 591 bp host repeat, were identical. Further inwards,

459  however, the order and arrangement of EVE fragments diverged.

460  As in the S97-EVE, the central region of the S35-EVE1 contained a block of 4 x 661 bp

461  RUs each comprised of the same inward facing (A) 398 bp and (B) 263 bp repeats but

462  ordered 5'|A-B|B-A|A-B|B-A|3', thus making a 2877 bp inverted repeat with an inversion

463  point at the RU2-RU3 boundary. Also, like the 97-EVE, each S35-EVE RU was flanked

464  by same 18 bp inverted repeat sequence, with those upstream of RU1 and downstream

465  of RU4 extending 17-33 bp into a 41 bp imperfect inverted repeat sequence located

466  immediately upstream toward the 5'-genome termini in some IHHNV strains

467  (Supplementary Figure 2). However, unlike the RU block in the S97-EVE, each of the

468  three internal S35-EVE RU boundaries comprised of 2 x 18 bp inverted repeats flanking

469  the complete 41 bp imperfect inverted repeat (Supplementary Figure 2). This revised

470  the RU junction to the inversion point in longer imperfect inverted repeat, rather than the

471  inversion point of the 18 bp inverted repeat. DNAfold and RNAfold analyses showed

472  that the 41 bp inverted repeat and its reverse compliment sequence could fold into

473  stable hairpin structures as either ssDNA ($\Delta G$ = -14.18/-14.86, Tm = 73.6/75.6℃) or

474  ssRNA ($\Delta G$ = -22.50/-25.00).

475  The larger S35-EVE2 sequence differed in the arrangement and homology of up to

476  eight RUs, possibly composed of two entirely duplicated inward-facing EVE fragments

477  (Table 3). The IHHNV-EVE fragments in S35-EVE1 contained 72% of the Au2005

478  IHHNV-EVE sequence with 98.8% homology, on average. In contrast, IHHNV-EVE

479  fragments in S35-EVE2 region only contained 53% of the IHHNV-EVE sequence with

480  97.5% homology, on average.

481  BLASTn searches of the genome assembly of a *P. monodon* from Vietnam

482  (Pmod26D_v1, GCA_007890405.1, Van Quyen et al. 2020), using the 9,045 bp S97-

483  EVE and 3,832 bp linear Au2005 Type A IHHNV-EVE sequences identified 3 short

484  contigs (*VIGR010059916.1, 4,003 nt; VIGR010168684.1, 2,220 nt; VIGR010211091.1,*

485  1,917 bp) also comprised of jumbled IHHNV-EVE Type A-like fragments (Table 3). In

486  two of the contigs, the stretches of jumbled EVE fragments neighbored either a

487  complete (590 bp) or incomplete (356 bp) host repeat sequences like those flanking the

488  S97-EVE. BLASTn searches of a genome assembly of a *P. monodon* from Shenzhen,

489  China (Pmon_WGS_v1, GCA_002291185.1) also identified evidence of an EVE

490  comprised of jumbled IHHNV genome fragments (Table 3), and despite contig lengths

491  being short, it was also being flanked by the same repeated host sequence flanking the

492  S97-EVE (data not shown). While more complete higher quality genome assemblies

493 would add confidence, the insertion locations of the jumbled EVEs present in the

494 genome assemblies of the *P. monodon* from Vietnam and China appear shared with

495 those the Australian S97-EVE and Thai S35-EVE1, with the second less-related

496 jumbled S35-EVE2 in the Thai genome residing at a nearby site. Interestingly, BLASTn

497 searches of the genome assemblies of *P. monodon* from Australia, Thailand, Vietnam,

498 or China identified no evidence of linear IHHNV-EVE forms.

**Origins and implications of jumbled IHHNV-EVEs:**

500 While varying in lengths, the amalgamations of reordered, inverted, and repeated

501 IHHNV genome fragments comprising the EVEs detected in Scaffold_97 (S97) of the

502 Australian *P. monodon* assembly (this study) and in Scaffold_35 (S35) of the Thai *P.*

503 *monodon* assembly (Uengwetwanit et al. 2020) share an integration site as well as

504 structural and sequence similarities with the partial EVE sequences detected in short

505 contigs of genome assemblies of *P. monodon* originating from Vietnam and China (as

506 outlined above). These similarities are suggestive of a progenitor IHHNV genome

507 becoming stably integrated as an EVE prior to *P. monodon* becoming dispersed widely

508 across its current distribution range. Such an ancient event would also support

509 differences noted for example in EVE fragment composition, central RU numbers, and

510 the nature of the conserved inverted-repeat sequences defining the boundaries of the

511 RUs. Furthermore, the conservation of the inverted-repeat sequences at the RU

512 boundaries and their potential to form stable ssDNA folding structures suggests a

513 potential role in their apparent multiplication.

514 The IHHNV P2 RNA transcriptional promoter motifs, including the 18 bp inverted repeat

515 sequences and TATA-box (Dhar et al. 2014; Shike et al. 2000; Silva et al. 2014), at the

516 RU boundaries have potential to facilitate transcription of various virus-specific sense

517 and antisense ssRNA sequences. RNA transcribed from them would then be capable of

518 forming long virus-specific dsRNA or hairpin dsRNAs, potentially in high abundance due

519 to their repeated nature. If so, such virus-specific antisense RNAs or dsRNA forms

520 processed through the RNA interference (RNAi) machinery of *P. monodon* (Attasart et

521 al. 2010; Attasart et al. 2011; Dhar et al. 2014; Su et al. 2008) could provide resilience

522 against IHHNV infections progressing to become acute and cause disease. Such an

523   advantage might promote the selection of *P. monodon* carrying this form of IHHNV-

524   EVE, particularly in circumstances when shrimp are specifically selected or bred for

525   aquaculture robustness. Selection for the EVE over several years would also be

526   consistent with the viral accommodation model hypothesized to involve farmed shrimp

527   acquiring and/or selected for an ability to mount elevated antisense ssRNA-based

528   and/or dsRNA-based anti-viral responses (Flegel 2007, 2020; Flegel 2009).

529   EVEs comprised of reordered, inverted, repeated and missing IHHNV genome

530   fragments would be expected to invalidate many PCR tests either designed specifically,

531   or found through use, to amplify IHHNV-EVE dsDNA sequences (Cowley et al. 2018;

532   Rai et al. 2009; Rai et al. 2012; Saksmerprome et al. 2011; Tang et al. 2007). As

533   examples, the 356 bp sequence targeted by the 77102F/77353R primer set (Nunan et

534   al. 2000) found to amplify both viral ssDNA and EVE dsDNA sequences existed in the

535   S97-EVE and S35-EVE1, but not in the S35-EVE2 sequence. However, nucleotide

536   mismatches at the 3' terminal position of both primers and at four other positions in the

537   18-mer 77353R primer would likely compromise the capacity of this primer set test to

538   amplify these EVEs. In contrast, neither EVE sequence possessed intact fragments

539   spanning regions amplified by primer sets 392F/R (392 bp) and 389F/R (389 bp)

540   recommended by the OIE as useful for amplifying divergent IHHNV strains as well as

541   IHHNV-EVE Type A and B sequences, or primer set MG831F/R (831 bp) designed

542   specifically to amplify known linear IHHNV-EVE types (Tang et al. 2007). Similarly, the

543   region targeted by a real-time PCR primer set designed to specifically amplify IHHNV-

544   EVE Type A sequences was absent from the S97-EVE and S35-EVE1, but present,

545   albeit with some primer mismatches, in the S35-EVE2 sequence (Cowley et al., 2018).

546   Variability among individual *P. monodon* in EVE sequences amplified by a suite of 10

547   PCR primer sets covering overlapping regions of complete linear IHHNV-EVE sequence

548   have been interpreted to suggest the random integration of IHHNV genome fragments

549   (Saksmerprome et al. 2011). While the jumbled fragments in the IHHNV-EVEs

550   described here might explain these, the diversity in EVE makeup suggested by these

551   data would require jumbled EVEs to be characterized in larger numbers of *P. monodon,*

552   or other penaeid species susceptible to IHHNV infection. Such broader information will

553　also be important to devising PCR methods to detect jumbled IHHNV-EVE sequences

554　more reliably.

**Conclusions:**

556　Using PacBio long-read data with Illumina short-read polishing together with 10X

557　Genomics and Hi-C scaffolding, this study generated a draft genome assembly and

558　annotation of a black tiger shrimp (*Penaeus monodon*) originating from Australia. The

559　assembly represents the first to be produced from this geographically isolated and

560　genetically distinct population (Vu et al. 2021). The assembly therefore adds to the

561　genetic resources available for *P. monodon* and Penaeid shrimp in general, and will

562　assist investigations into their evolution and genome expansion resulting from

563　transposable elements. Of the *P. monodon* genome features, the high prevalence of

564　general repeats is the most remarkable, and especially the high content of SSRs even

565　in comparison to other crustacean species. Another unexpected feature was the

566　existence of a previously undescribed IHHNV endogenous viral element (EVE) located

567　between a repeated host sequence. Rather than being comprising of a linear sequence

568　of all or part of the ~3.9 kb IHHNV genome, the EVE comprised of a conglomerate of

569　reordered, inverted, and repeated IHHNV genome fragments. Searches of genome

570　assemblies available for *P. monodon* from Thailand, Vietnam and China indicated with

571　variable confidence, depending on assembly quality, that each contained a similarly

572　jumbled IHHNV-EVE inserted at the same genome location. The fragmented and

573　rearranged nature of these EVEs has implications for detecting them with currently

574　available PCR tests. The presence of multiple inverted sequences including multiple

575　IHHNV RNA transcription promoter elements also has implications for them expressing

576　virus-specific dsRNA capable of interfering with exogenous IHHNV replication. The

577　complexity of the rearranged IHHNV genome fragments comprising the EVEs begs

578　many questions related to how long they have existed in the genomes of genetically

579　diverse *P. monodon*, as well as to what processes have led to their integration at a

580　specific genome location, to the IHHNV genome fragments becoming rearranged and to

581　the apparent multiplication of a repeat unit comprised of highly defined inverted

582　sequences derived from the 5'-terminal region of the IHHNV genome.

597

598 **Table 1** Illumina, PacBio, 10X Genomics, and DoveTail sequencing data used for the

599 assembly and scaffolding of the black tiger shrimp genome.

| Sequencing Platform | Paired End Reads | Yield (Gb) | Coverage | GenBank accessions |
|---|---|---|---|---|
| Illumina (250 bp PE) | 315 M | 158 | 72 X | SRR10713996, SRR10713997 |
| PacBio Sequel | N/A | 165 | 75 X | SRR10713990 - SRR10713995 SRR10713998 - SRR10714025 |
| 10X Genomics (250 bp PE) | 987 M | 494 | 224 X | N/A |
| DoveTail (100 bp PE) | 1.2 B | 119 | 54 X | N/A |

600

601 **Table 2** Summary of assembly statistics for the Australian and Thai *P. monodon*, and *P.*

602 *vannamei* genomes.

| Metrics | *P. monodon* (Australia) | *P. monodon* (Thailand) | *P. vannamei* |
|---|---|---|---|
| # contigs | 47,607 | 70,380 | 50,304 |
| Largest contig | 1,147,530 | 1,387,722 | 739,419 |
| Total length of contigs | 1.89 Gb | 2.39 Gb | 1.62 Gb |
| Contig N50 | 78 kb | 79 kb | 58 kb |
| # Scaffolds | 31,922 | 44 | - |
| Largest scaffold | 21.70 Mb | 65.87 Mb | - |
| Total length of scaffolds | 1.89 Gb | 1.99 Gb | 1.66 Gb |
| Scaffold N50 | 0.50 Mb | 49.0 Mb | 0.60 Mb |
| Projected Genome Size | 2.20 Gb | 2.20 Gb | 2.45 Gb |
| Percentage Covered By Scaffolds | 86.1% | 90.3% | 67.7% |
| GC (%) | 35.6 | 36.6 | 35.7 |
| Complete BUSCOs (C) | 86.8 | 87.9 | 78.0 |
| Complete and single-copy BUSCOs | 85.8 | 84.8 | 74.0 |
| Complete and duplicated BUSCOs | 1.0 | 3.1 | 4.0 |
| Fragmented BUSCOs (F) | 4.5 | 4.0 | 4.0 |
| Missing BUSCOs (M) | 8.7 | 8.0 | 18.0 |
| No. predicted gene models | 35,517 | 31,640 | 25,596 |
| No. of protein coding genes | 25,809 | 30,038 | - |
| No. genes annotated in | 17,158 | 20,615 | - |
| References | This study | Uengwetwanit et al. (2020) | Zhang et al. (2019) |

603

604

605 **Table 3** Detection and notable features of IHHNV-EVE sequences identified in other

606 genomes of *P. monodon*.

| Reference Genome IDs | Notable EVE features | | | | |
|---|---|---|---|---|---|
| | Start | End | Length (bp) | Orientation | Homology (%) |
| *P. monodon* Thailand (Uengwetwanit et al. 2020) | | | | | |
| *Scaffold 35 EVE-1* | 770,236 | 778,124 | 7,888 | | |
| RU1 | 772,730 | 773,391 | 661 | minus | 99.9 |
| RU2 | 773,450 | 774,111 | 661 | plus | 100.0 |
| RU3 | 774,170 | 774,831 | 661 | minus | 99.9 |
| RU4 | 774,890 | 775,551 | 661 | plus | 97.9 |
| *Scaffold 35 EVE-2* | 862,618 | 878,928 | 16,310 | | |
| RU1 | 866,534 | 867,145 | 611 | minus | 79.4 |
| RU2 | 867,204 | 867,791 | 587 | plus | 81.3 |
| RU3 | 867,840 | 868,467 | 627 | minus | 83.5 |
| RU4 | 868,515 | 869,130 | 615 | plus | 80.0 |
| RU5 | 872,127 | 872,754 | 627 | plus | 78.9 |
| RU6 | 872,799 | 873,434 | 635 | minus | 90.0 |
| RU7 | 873,492 | 874,152 | 660 | plus | 97.2 |
| RU8 | 875,469 | 876,168 | 699 | plus | 92.1 |
| *P. monodon* Vietnam (Pmod26D_v1; GCA_007890405.1) | | | | | |
| VIGR010059916.1 EVE (4,003 bp) | | | 4,003 | | 98.4 |
| VIGR010211091.1 EVE (1,917 bp) | | | 1,917 | | 99.0 |
| VIGR010168684.1 EVE (2,220 bp) | | | 2,220 | | 98.9 |
| *P. monodon* China (Pmon_WGS_v1, GCA_002291185.1) | | | | | |
| gb|NIUS011382605.1 (645 bp) | | | 645 | | 98.9 |
| gb|NIUS011109800.1 (848 bp) | | | 848 | | 98.3 |

607

608

609     ***References***

610     Angthong, P., T. Uengwetwanit, W. Pootakham, K. Sittikankaew, C. Sonthirod *et al.*,
611          2020 Optimization of high molecular weight DNA extraction methods in shrimp
612          for a long-read sequencing platform. *PeerJ* 8:e10340.

613     Attasart, P., R. Kaewkhaw, C. Chimwai, U. Kongphom, O. Namramoon *et al.*, 2010
614          Inhibition of Penaeus monodon densovirus replication in shrimp by double-
615          stranded RNA. *Archives of virology* 155 (6):825-832.

616     Attasart, P., R. Kaewkhaw, C. Chimwai, U. Kongphom, and S. Panyim, 2011 Clearance
617          of Penaeus monodon densovirus in naturally pre-infected shrimp by combined
618          ns1 and vp dsRNAs. *Virus research* 159 (1):79-82.

619     Benson, G., 1999 Tandem repeats finder: a program to analyze DNA sequences.
620          *Nucleic acids research* 27 (2):573-580.

621     Bikard, D., S. Julie-Galau, G. Cambray, and D. Mazel, 2010 The synthetic integron: an
622          in vivo genetic shuffling device. *Nucleic acids research* 38 (15):e153-e153.

623     Campbell, M.S., C. Holt, B. Moore, and M. Yandell, 2014 Genome annotation and
624          curation using MAKER and MAKER-P. *Current protocols in bioinformatics* 48
625          (1):4.11. 11-14.11. 39.

626     Cantarel, B.L., I. Korf, S.M. Robb, G. Parra, E. Ross *et al.*, 2008 MAKER: an easy-to-
627          use annotation pipeline designed for emerging model organism genomes.
628          *Genome research* 18 (1):188-196.

629     Chandhini, S., and V.J. Rejish Kumar, 2019 Transcriptomics in aquaculture: current
630          status and applications. *Reviews in Aquaculture* 11 (4):1379-1397.

631     Cooke, I.R., H. Ying, S. Foret, P. Bongaerts, J.M. Strugnell *et al.*, 2020 Signatures of
632          selection in the coral holobiont reveal complex adaptations to inshore
633          environments driven by Holocene climate change. *bioRxiv*.

634     Cotmore, S., and P. Tattersall, 1996 Parvovirus DNA replication. *DNA replication in*
635          *eukaryotic cells. Cold Spring Harbor Laboratory Press, Cold Spring Harbor,*
636          *NY*:799-813.

637     Cowley, J.A., M. Rao, and G.J. Coman, 2018 Real-time PCR tests to specifically detect
638          IHHNV lineages and an IHHNV EVE integrated in the genome of Penaeus
639          monodon. *Diseases of aquatic organisms* 129 (2):145-158.

640      Dhar, A.K., K.N. Kaizer, Y.M. Betz, T.N. Harvey, and D.K. Lakshman, 2011
641          Identification of the core sequence elements in Penaeus stylirostris densovirus
642          promoters. *Virus genes* 43 (3):367-375.
643      Dhar, A.K., K.N. Kaizer, and D.K. Lakshman, 2010 Transcriptional analysis of Penaeus
644          stylirostris densovirus genes. *Virology* 402 (1):112-120.
645      Dhar, A.K., D.K. Lakshman, S. Natarajan, F.T. Allnutt, and N.A. van Beek, 2007
646          Functional characterization of putative promoter elements from infectious
647          hypodermal and hematopoietic necrosis virus (IHHNV) in shrimp and in insect
648          and fish cell lines. *Virus research* 127 (1):1-8.
649      Dhar, A.K., R. Robles-Sikisaka, V. Saksmerprome, and D.K. Lakshman, 2014 Biology,
650          genome organization, and evolution of parvoviruses in marine shrimp. *Advances*
651          *in virus research* 89:85-139.
652      Dobin, A., C.A. Davis, F. Schlesinger, J. Drenkow, C. Zaleski *et al.*, 2013 STAR:
653          ultrafast universal RNA-seq aligner. *Bioinformatics* 29 (1):15-21.
654      FAO, 2020 *The state of world fisheries and aquaculture 2020: Sustainability in action*:
655          Food and Agriculture Organization of the United Nations.
656      Fisher, S., A. Barry, J. Abreu, B. Minie, J. Nolan *et al.*, 2011 A scalable, fully automated
657          process for construction of sequence-ready human exome targeted capture
658          libraries. *Genome Biol* 12 (1):R1.
659      Flegel, T., 2007 Update on viral accommodation, a model for host-viral interaction in
660          shrimp and other arthropods. *Developmental & Comparative Immunology* 31
661          (3):217-231.
662      Flegel, T., 2020 Research progress on viral accommodation 2009 to 2019.
663          *Developmental & Comparative Immunology*:103771.
664      Flegel, T.W., 2009 Hypothesis for heritable, anti-viral immunity in crustaceans and
665          insects. *Biology Direct* 4 (1):1-8.
666      Ghurye, J., M. Pop, S. Koren, D. Bickhart, and C.-S. Chin, 2017 Scaffolding of long read
667          assemblies using long range contact information. *BMC genomics* 18 (1):1-11.
668      Guppy, J.L., D.B. Jones, S.R. Kjeldsen, A. Le Port, M.S. Khatkar *et al.*, 2020
669          Development and validation of a RAD-Seq target-capture based genotyping

670     assay for routine application in advanced black tiger shrimp (Penaeus monodon)

671     breeding programs.

672  Hauton, C., 2017 Recent progress toward the identification of anti-viral immune

673     mechanisms in decapod crustaceans. *Journal of invertebrate pathology* 147:111-

674     117.

675  Hollenbeck, C.M., and I.A. Johnston, 2018 Genomic tools and selective breeding in

676     molluscs. *Frontiers in genetics* 9:253.

677  Holt, C., and M. Yandell, 2011 MAKER2: an annotation pipeline and genome-database

678     management tool for second-generation genome projects. *BMC bioinformatics*

679     12 (1):491.

680  Houston, R.D., T.P. Bean, D.J. Macqueen, M.K. Gundappa, Y.H. Jin *et al.*, 2020

681     Harnessing genomics to fast-track genetic improvement in aquaculture. *Nature*

682     *Reviews Genetics* 21 (7):389-409.

683  Huang, S.-W., Y.-Y. Lin, E.-M. You, T.-T. Liu, H.-Y. Shu *et al.*, 2011 Fosmid library end

684     sequencing reveals a rarely known genome structure of marine shrimp Penaeus

685     monodon. *BMC genomics* 12 (1):242.

686  Huerlimann, R., N.M. Wade, L. Gordon, J.D. Montenegro, J. Goodall *et al.*, 2018 De

687     novo assembly, characterization, functional annotation and expression patterns

688     of the black tiger shrimp (Penaeus monodon) transcriptome. *Scientific reports* 8

689     (1):1-14.

690  Kawato, S., K. Nishitsuji, A. Arimoto, K. Hisata, M. Kawamitsu *et al.*, 2021 Genome and

691     transcriptome assemblies of the kuruma shrimp, Marsupenaeus japonicus. *G3*

692     *Genes| Genomes| Genetics*.

693  Kim, D., J.M. Paggi, C. Park, C. Bennett, and S.L. Salzberg, 2019 Graph-based

694     genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature*

695     *biotechnology* 37 (8):907-915.

696  Koressaar, T., and M. Remm, 2007 Enhancements and modifications of primer design

697     program Primer3. *Bioinformatics* 23 (10):1289-1291.

698  Krabsetsve, K., B.R. Cullen, and L. Owens, 2004 Rediscovery of the Australian strain of

699     infectious hypodermal and haematopoietic necrosis virus. *Diseases of aquatic*

700     *organisms* 61 (1-2):153-158.

701   Kulkarni, A., S. Krishnan, D. Anand, S. Kokkattunivarthil Uthaman, S.K. Otta *et al.*, 2021
702           Immune responses and immunoprotection in crustaceans with special reference
703           to shrimp. *Reviews in Aquaculture* 13 (1):431-459.

704   Kurtzer, G.M., V. Sochat, and M.W. Bauer, 2017 Singularity: Scientific containers for
705           mobility of compute. *PloS one* 12 (5).

706   Kotin, R.M., R.M. Linden, and K.I. Berns, 1992 Characterization of a preferred site on
707           human chromosome 19q for integration of adeno-associated virus DNA by non-
708           homologous recombination. *The EMBO journal* 11 (13): 5071-5078.

709   Linden, R.M., E. Winocour, and K.I. Berns (1996). The recombination signals for adeno-
710           associated virus site-specific integration. *PNAS USA* 93 (15): 7966-7972.

711   Nunan, L.M., B.T. Poulos, and D.V. Lightner, 2000 Use of Polymerase Chain Reaction
712           for the Detection of Infectious Hypodermal and Hematopoietic Necrosis Virus in
713           Penaeid Shrimp. *Marine Biotechnology* 2 (4):319-328.

714   Pertea, M., G.M. Pertea, C.M. Antonescu, T.-C. Chang, J.T. Mendell *et al.*, 2015
715           StringTie enables improved reconstruction of a transcriptome from RNA-seq
716           reads. *Nature biotechnology* 33 (3):290.

717   Rai, P., B. Pradeep, I. Karunasagar, and I. Karunasagar, 2009 Detection of viruses in
718           Penaeus monodon from India showing signs of slow growth syndrome.
719           *Aquaculture* 289 (3):231-235.

720   Rai, P., M.P. Safeena, K. Krabsetsve, K. La Fauce, L. Owens *et al.*, 2012 Genomics,
721           Molecular Epidemiology and Diagnostics of Infectious hypodermal and
722           hematopoietic necrosis virus. *Indian J Virol* 23 (2):203-214.

723   Ruan, J., and H. Li, 2019 Fast and accurate long-read assembly with wtdbg2. *Nature*
724           *Methods*:1-4.

725   Saksmerprome, V., S. Jitrakorn, K. Chayaburakul, S. Laiphrom, K. Boonsua *et al.*, 2011
726           Additional random, single to multiple genome fragments of Penaeus stylirostris
727           densovirus in the giant tiger shrimp genome have implications for viral disease
728           diagnosis. *Virus research* 160 (1-2):180-190.

729   Schnepp, B C., R.L. Jensen, C.L. Chen, P.R. Johnson and K.R. Clark, 2005
730           Characterization of adeno-associated virus genomes isolated from human
731           tissues. *Journal of virology* 79(23): 14793-14803.

732 Sellars, M.J., L. Dierens, S. McWilliam, B. Little, B. Murphy *et al.*, 2014 Comparison of
733       microsatellite and SNP DNA markers for pedigree assignment in Black Tiger
734       shrimp, *Penaeus monodon. Aquaculture Research* 45 (3):417-426.

735 Servant, N., N. Varoquaux, B.R. Lajoie, E. Viara, C.-J. Chen *et al.*, 2015 HiC-Pro: an
736       optimized and flexible pipeline for Hi-C data processing. *Genome biology* 16
737       (1):1-11.

738 Shike, H., A.K. Dhar, J.C. Burns, C. Shimizu, F.X. Jousset *et al.*, 2000 Infectious
739       hypodermal and hematopoietic necrosis virus of shrimp is related to mosquito
740       brevidensoviruses. *Virology* 277 (1):167-177.

741 Silva, D.C., A.R. Nunes, D.I. Teixeira, J.P.M. Lima, and D.C. Lanza, 2014 Infectious
742       hypodermal and hematopoietic necrosis virus from Brazil: Sequencing,
743       comparative analysis and PCR detection. *Virus research* 189:136-146.

744 Simão, F.A., R.M. Waterhouse, P. Ioannidis, E.V. Kriventseva, and E.M. Zdobnov, 2015
745       BUSCO: assessing genome assembly and annotation completeness with single-
746       copy orthologs. *Bioinformatics* 31 (19):3210-3212.

747 Su, J., D.T. Oanh, R.E. Lyons, L. Leeton, M.C. van Hulten *et al.*, 2008 A key gene of the
748       RNA interference pathway in the black tiger shrimp, *Penaeus monodon*:
749       identification and functional characterisation of Dicer-1. *Fish & shellfish*
750       *immunology* 24 (2):223-233.

751 Taengchaiyaphum, S., P. Buathongkam, S. Sukthaworn, P. Wongkhaluang, K.
752       Sritunyalucksana *et al.*, 2021 Shrimp parvovirus circular DNA fragments arise
753       from both endogenous viral elements (EVE) and the infecting virus. *bioRxiv*.

754 Tang, B., D. Zhang, H. Li, S. Jiang, H. Zhang *et al.*, 2020 Chromosome-level genome
755       assembly reveals the unique genome evolution of the swimming crab (*Portunus*
756       *trituberculatus*). *GigaScience* 9 (1):giz161.

757 Tang, K.F., S.A. Navarro, and D.V. Lightner, 2007 PCR assay for discriminating
758       between infectious hypodermal and hematopoietic necrosis virus (IHHNV) and
759       virus-related sequences in the genome of *Penaeus monodon. Diseases of*
760       *aquatic organisms* 74 (2):165-170.

761   Uengwetwanit, T., W. Pootakham, I. Nookaew, C. Sonthirod, P. Angthong *et al.*, 2020 A
762        chromosome-level assembly of the black tiger shrimp (*Penaeus monodon*)
763        genome facilitates the identification of novel growth-associated genes. *bioRxiv*.
764   Untergasser, A., I. Cutcutache, T. Koressaar, J. Ye, B.C. Faircloth *et al.*, 2012
765        Primer3—new capabilities and interfaces. *Nucleic acids research* 40 (15):e115-
766        e115.
767   Van Quyen, D., H.M. Gan, Y.P. Lee, D.D. Nguyen, T.H. Nguyen *et al.*, 2020 Improved
768        genomic resources for the black tiger prawn (*Penaeus monodon*). *Marine*
769        *Genomics*:100751.
770   Vu, N.T., K.R. Zenger, C.N. Silva, J.L. Guppy, and D.R. Jerry, 2021 Population
771        structure, genetic connectivity, and signatures of local adaptation of the Giant
772        Black tiger shrimp (*Penaeus monodon*) throughout the Indo-Pacific region.
773        *Genome biology and evolution* 13 (10):evab214.
774   Walker, B.J., T. Abeel, T. Shea, M. Priest, A. Abouelliel *et al.*, 2014 Pilon: an integrated
775        tool for comprehensive microbial variant detection and genome assembly
776        improvement. *PloS one* 9 (11).
777   Weisenfeld, N.I., V. Kumar, P. Shah, D.M. Church, and D.B. Jaffe, 2017 Direct
778        determination of diploid genome sequences. *Genome research* 27 (5):757-767.
779   Weisenfeld, N.I., S. Yin, T. Sharpe, B. Lau, R. Hegarty *et al.*, 2014 Comprehensive
780        variation discovery in single human genomes. *Nature genetics* 46 (12):1350.
781   Yang, W., N.T. Tran, C.-H. Zhu, D.-F. Yao, J.J. Aweya *et al.*, 2021 Immune priming in
782        shellfish: A review and an updating mechanistic insight focused on cellular and
783        humoral responses. *Aquaculture* 530:735831.
784   Yang, C.C., X. Xiao, X. Zhu, D.C. Ansardi, N.D. Epstein *et al.*, 1997 Cellular
785        recombination pathways and viral terminal repeat hairpin structures are sufficient
786        for adeno-associated virus integration *in vivo* and *in vitro. Journal of virology* 71
787        (12): 9231-9247.
788   Yeo, S., L. Coombe, J. Chu, R.L. Warren, and I. Birol, 2017 ARCS: assembly roundup
789        by chromium scaffolding. *bioRxiv*:100750.

790 Yu, Y., X. Zhang, J. Yuan, F. Li, X. Chen *et al.*, 2015 Genome survey and high-density
791         genetic map construction provide genomic and genetic resources for the Pacific
792         White Shrimp *Litopenaeus vannamei*. *Scientific reports* 5 (1):1-14.
793 Yuan, J., X. Zhang, F. Li, and J. Xiang, 2021a Genome sequencing and assembly
794         strategies and a comparative analysis of the genomic characteristics in penaeid
795         shrimp species. *Frontiers in genetics* 12.
796 Yuan, J., X. Zhang, C. Liu, Y. Yu, J. Wei *et al.*, 2018 Genomic resources and
797         comparative analyses of two economical penaeid shrimp species, *Marsupenaeus*
798         *japonicus* and *Penaeus monodon*. *Marine Genomics* 39:22-25.
799 Yuan, J., X. Zhang, M. Wang, Y. Sun, C. Liu *et al.*, 2021b Simple sequence repeats
800         drive genome plasticity and promote adaptive evolution in penaeid shrimp.
801         *Communications biology* 4 (1):1-14.
802 Yue, G., and L. Wang, 2017 Current status of genome sequencing and its applications
803         in aquaculture. *Aquaculture* 468:337-347.
804 Zdobnov, E.M., F. Tegenfeldt, D. Kuznetsov, R.M. Waterhouse, F.A. Simao *et al.*, 2017
805         OrthoDB v9. 1: cataloging evolutionary and functional annotations for animal,
806         fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic acids research* 45
807         (D1):D744-D749.
808 Zenger, K., M. Khatkar, D. Jerry, and H. Raadsma, 2017 The next wave in selective
809         breeding: implementing genomic selection in aquaculture, pp. 105-112 in *Proc.*
810         *Assoc. Advmt. Anim. Breed. Genet*.
811 Zhang, X., J. Yuan, Y. Sun, S. Li, Y. Gao *et al.*, 2019 Penaeid shrimp genome provides
812         insights into benthic adaptation and frequent molting. *Nature communications* 10
813         (1):1-14.
814
815