# *De novo* assembly, characterization, functional annotation and expression patterns of the black tiger shrimp (*Penaeus monodon*) transcriptome

Roger Huerlimann[1,2*†], Nicholas M Wade[1,3†], Lavinia Gordon[1,4], Juan D Montenegro[4], Jake Goodall[1,3], Sean McWilliam[3], Matthew Tinning[1,4], Kirby Siemering[1,4], Erika Giardina[1,5], Dallas Donovan[1,5], Melony J Sellars[1,3], Jeff A Cowley[1,3], Kelly Condon[1,2], Greg J Coman[1,7], Mehar S Khatkar[1,6], Herman W Raadsma[1,6], Gregory Maes[8,9], Kyall R Zenger[1,2], and Dean R Jerry[1,2].

*Corresponding author: roger.huerlimann@jcu.edu.au

[†]These two authors contributed equally to the writing of this publication

Affiliations:

1. ARC Research Hub for Advanced Prawn Breeding

2. Centre for Sustainable Tropical Fisheries and Aquaculture, College of Science and Engineering, James Cook University, Townsville, QLD 4811, Australia

3. CSIRO Agriculture and Food, Integrated Sustainable Aquaculture Program, 306 Carmody Road, St Lucia, QLD 4067 Australia

4. Australian Genome Research Facility Ltd, The Walter and Eliza Hall Institute, 1G Royal Parade, Parkville VIC 3050

5. Seafarms Group Ltd, Level 11 225 St Georges Terrace Perth, WA 6000.

6. Sydney School of Veterinary Science, Faculty of Science, The University of Sydney, NSW, Australia

7. CSIRO Agriculture and Food, Integrated Sustainable Aquaculture Program, 144 North Street, Woorim, QLD 4507 Australia

8. Laboratory of Biodiversity and Evolutionary Genomics, KU Leuven, Leuven, 3000, Belgium

9. Center for Human Genetics, UZ Leuven- Genomics Core, KU Leuven, Leuven, 3000, Belgium

## Declarations

Availability of data and material

Raw data, assembly and bioinformatics scripts will be made freely available online upon publication.

Funding

This work was funded by the Australian Research Council (ARC) Industrial Transformation Research Hub scheme, awarded to James Cook University and in collaboration with the Commonwealth Scientific Industrial Research Organisation (CSIRO), the Australian Genome Research Facility (AGRF), the University of Sydney and Seafarms Group Pty Ltd.

Authors' contributions

RH: conceptualised, developed and oversaw the project, performed sampling and carried out RNA extractions, developed and performed the transcriptome assembly, quality assessment and differential gene analysis, and wrote the manuscript. NMW: conceptualised and developed the project, performed sampling, led components of data analysis and interpretation and wrote the manuscript. LG: developed the transcriptome assembly bioinformatics pipeline. JDM: carried out the lncRNA analysis and assisted with the transcriptome assembly bioinformatics pipeline. JG: carried out sampling and RNA extractions, and assisted with development of the transcriptome assembly bioinformatics pipeline, and reviewed the manuscript. SM: assisted with the bioinformatic analysis of the differential gene expression data and reviewed manuscript. MT: oversaw the library preparation and sequencing, and reviewed the manuscript. KS: conceptualised and developed the project and reviewed manuscript. EG: reared and sampled the larval stages. DD: conceptualised and developed the project. Coordinated facilities and resources for larval and adult prawn production and reviewed manuscript. MS: provided prawn tissue samples, conceptualised and developed the project, assisted with data interpretation and reviewed manuscript. JC: assisted with the interpretation and writing of the viral analysis, and edited manuscript. KC: assisted with the interpretation and writing of the viral analysis, and reviewed manuscript. GC: conceptualised and developed the project and reviewed manuscript. MK: conceptualised and developed the project and reviewed manuscript. HR: conceptualised and developed the project and reviewed manuscript. GM: conceptualised and developed the project, performed sampling, provided advice on sequencing strategies and data interpretation and reviewed manuscript. KRZ: conceptualised and developed the project and reviewed manuscript. DRJ: conceptualised and developed the project, oversaw coordination of project activities and reviewed manuscript

All authors read and approved the final manuscript.

1 ## **<u>Abstract</u>**

2 The black tiger shrimp (*Penaeus monodon*) remains the second most widely cultured

3 shrimp species globally. However, issues with disease and domestication have seen

4 production levels stagnate over the past two decades. To help identify innovative

5 solutions needed to resolve bottlenecks hampering the culture of this species, it is

6 important to generate genetic and genomic resources. Towards this aim, we have

7 produced the most complete publicly available *P. monodon* transcriptome database

8 to date. The assembly was carried out in multiple assemblers using 2x125 bp HiSeq

9 data from PolyA selected, ribo-depleted RNA extracted from nine adult tissues and

10 eight early life-history stages. In total, approximately 700 million high-quality

11 sequence reads were obtained and assembled into 236,388 clusters. These were

12 then further segregated into 99,203 adult tissue specific clusters, and 58,678 early

13 life-history stage specific clusters. The final transcriptome had a high TransRate

14 score of 0.37, with 88% of all reads successfully mapping back to the transcriptome.

15 BUSCO statistics showed the assembly to be highly complete with low

16 fragmentation, few genes missing, but higher redundancy or transcript duplication

17 (Complete: 98.2% (Duplicated: 51.3%), Fragmented: 0.8%, Missing: 1.0%), and to

18 greatly exceed the completeness of existing *P. monodon* transcriptomes. While

19 annotation rates were low (approximately 30%), as is typical for a non-model

20 organisms, annotated transcript clusters were successfully mapped to several

21 hundred functional KEGG pathways. To help address the lack of annotation,

22 transcripts were clustered into groups within tissues and early life-history stages,

23 providing initial evidence for their roles in specific tissue functions, or developmental

24 transitions. Additionally, transcripts of shrimp viruses previously not known to occur

25 in Australia were also discovered. We expect the transcriptome to provide an

26 essential resource to investigate the molecular basis of commercially relevant-

27 significant traits in *P. monodon* and other shrimp species.

28

## Introduction

30 The black tiger shrimp *Penaeus monodon* belongs to the family Penaeidae and is
31 the second most widely farmed shrimp species globally[1]. However, disease and
32 limited progress in domestication and selective breeding of *P. monodon* continue to
33 hamper further expansion of the industry[2]. Modern genomic technologies have
34 significant potential to advance selective breeding programs; however, they require
35 complete, well annotated tissue-specific transcriptomic and genomic datasets. In
36 addition to assisting in genome assembly and creating linkage maps[3], a complete
37 transcriptome provides a potential resource for differential gene-expression
38 studies[4]), genome annotation[5], single nucleotide polymorphism discovery[6] and
39 genome scaffolding[7].

40 While genomic resources for Penaeid shrimp are increasing, they remain limited for
41 many species, including *P. monodon*. Previous research has focussed on
42 hepatopancreas, ovary, heart, muscle and eyestalk tissues[8,9], in male and female
43 gonads[10], and in response to infection with *Vibrio* bacterial species capable of
44 inducing acute hepatopancreatic necrosis disease[11]. In addition to such differential
45 gene-expression studies, genomic data from next generation sequencing (NGS)
46 methods has expanded greatly in recent years, particularly in the study of Pacific
47 white shrimp (*Litopenaeus vannamei*)[3,6,12-23]. Moreover, a transcriptome based on
48 eight tissues was assembled for the less well studied banana shrimp
49 *Fenneropenaeus merguiensis*[24], and genes involved in early embryonic specification
50 have been studied in *Marsupenaeus japonicus*[25]. Transcriptomics has also been
51 applied to *Penaeus merguiensis*[26-28] and the Chinese white shrimp *Fenneropenaeus*
52 *chinensis*[29,30] to investigate aspects of tissue-specific expression, stress tolerance
53 and viral infection. Despite these advances, a comprehensive transcriptome from
54 diverse tissue types and early life-history stages of *P. monodon* remains unavailable.

55 In order to address this deficiency, we report a highly complete transcriptome for *P.*
56 *monodon* that can be used as a broad basis for future genomics research. To this
57 effect, we sequenced three replicates each from nine different tissues types
58 (eyestalk, stomach, female gonad, male gonad, gill, haemolymph, hepatopancreas,
59 lymphoid organ and tail muscle) and one pooled replicate each from four larval
60 stages (embryo, nauplii, zoea, and mysis) and four post-larval stages ranging from
61 days 1, 4, 10 and 15. Additionally, transcript expression profiles unique to each type
62 and stage were determined, as well as identifying putative long non-coding RNA and
63 transcripts originating from viruses.

64

## **Material and Methods**

### *Sample taking and RNA extraction*

Tissues of *P. monodon* broodstock were collected from multiple individuals, immediately snap frozen on dry ice, and stored at -80°C until extraction (Table 1). All tissues except lymphoid organs were collected from wild broodstock caught off coastal waters near the border between the Northern Territory and Western Australia provided, which were provided by a commercial hatchery at Flying Fish Point, North Queensland, Australia. Lymphoid organ tissue was collected from wild prawns caught off the East Coast of Queensland. Larval and post-larval stages were collected from the same hatchery in pools of approximately 400 individuals per life stage, after four hours of starvation, and preserved in RNAlater (Thermo Fisher Scientific). All tissues and early life-history stages were sub-sampled in an RNase-free laboratory and total RNA was extracted using an RNeasy Universal extraction kit (QIAGEN) following manufacturer's instructions. RNA quantity and quality was estimated using a Nanodrop UV spectrophotometer (Thermo Fisher Scientific), and purity was further assessed using an Agilent Bioanalyzer (Agilent Technologies). RNA was selected from individual sample replicates based on Nanodrop spectra, RNA concentration, and Agilent Bioanalyzer traces (Table 1), in preference to using comparative tissues from the same individuals.

### *Illumina library preparation and sequencing*

Library preparation and sequencing was carried out at the Australian Genome Research Facility (AGRF). Upon arrival at the sequencing facility, the quality of the samples was checked using a Bioanalyzer RNA 6000 nano reagent kit (Agilent) and libraries were prepared using the TruSeq Stranded mRNA Library Preparation Kit (Illumina) according to established protocols. Final libraries were again checked using Tapestation DNA 1000 TapeScreen Assay (Agilent). Cluster generation was performed on a cBot with HiSeq PE Cluster Kit v4 - cBot and sequencing was done on a HiSeq 2500 using a HiSeq SBS Kit. The Hiseq 2500 was operating with HiSeq Control Software v2.2.68 and base-calling was performed with RTA v1.18.66.3. Samples in the second sequencing run were pooled and split across two lanes to reduce sequencing bias (Table 1).

97 **Table 1 | List of shrimp tissue types and early life-history stages**
98 **used for transcriptome sequencing.** PL = post-larval stages 1
99 (PL1), 4 (PL4), 10 (PL10), 15 (PL15)

| Shrimp ID | Sex | Tissue | Number of paired-end reads |
|---|---|---|---|
| PM_F_08 | Female | Eyestalk | 18,984,152 |
| | | Gill | 19,971,115 |
| | | Hepatopancreas | 18,831,682 |
| PM_F_02 | Female | Female Gonad | 21,338,933 |
| | | Haemolymph | 20,105,399 |
| | | Muscle | 20,361,299 |
| | | Stomach | 13,470,106 |
| PM_F_04 | Female | Female Gonad | 20,255,448 |
| | | Gill | 21,362,076 |
| | | Haemolymph | 20,247,206 |
| | | Stomach | 21,461,589 |
| PM_F_03 | Female | Female Gonad | 20,759,890 |
| PM_M_02 | Male | Eyestalk | 21,076,111 |
| | | Hepatopancreas | 19,029,973 |
| | | Male Gonad | 20,669,419 |
| | | Muscle | 20,129,858 |
| PM_M_04 | Male | Eyestalk | 22,250,295 |
| | | Gill | 20,396,956 |
| | | Haemolymph | 21,637,767 |
| | | Hepatopancreas | 20,854,492 |
| | | Male Gonad | 20,600,256 |
| | | Muscle | 22,464,431 |
| | | Stomach | 16,444,377 |
| PM_M_06 | Male | Male Gonad | 19,800,274 |
| PM_M_C2 | Male | Lymphoid Organ | 19,873,753 |
| PM_M_C3 | Male | Lymphoid Organ | 20,480,178 |
| PM_F_C1 | Female | Lymphoid Organ | 20,372,862 |
| Pool_E | | Embryo | 19,745,313 |
| Pool_N | | Nauplii | 18,310,089 |
| Pool_Z | | Zoea | 19,528,689 |
| Pool_M | | Mysis | 19,744,563 |
| Pool_PL1 | | PL1 | 19,815,103 |
| Pool_PL4 | | PL4 | 18,680,555 |
| Pool_PL10 | | PL10 | 18,773,667 |
| Pool_PL15 | | PL15 | 19,661,826 |

100

6

101 ### *Sequence quality control, assembly and annotation*

102 Raw sequence data was quality checked using FastQC[31] v0.11.5, and assembled

103 loosely following the Oyster River Protocol for Transcriptome Assembly[32]. In brief, all

104 sequences were collectively error-corrected using RCorrector[33] V3. Samples were

105 then assembled in Trinity[34] V2.3.2; grouped by individual shrimps, i.e. all tissues

106 from a specific shrimp were assembled together. Reads were trimmed harshly for

107 adapters and softly for Phred score <2 using Trimmomatic[35] V0.32; and then

108 normalized *in silico* within Trinity. The normalized forward and reverse reads

109 produced by Trinity were then used in BinPacker[36] V1.0, IDBA-Tran[37] V 1.1.1 using

110 K20, K30, K40, K50 and K60; and Bridger[38] version 2014-12-01. All resulting

111 transcriptomes were concatenated and merged using Evidential Gene[39], followed by

112 clustering using Transfuse V0.5.0 (https://github.com/cboursnell/transfuse) using a

113 similarity value of 0.98. Lastly, contigs <300 bp were removed to produce the final

114 transcriptome. The quality of the final assembly was assessed using TransRate[40]

115 V1.0.1, and BUSCO[41] V2 using the arthropoda_odb9 database[42]. Sequences were

116 annotated in Blast2Go[43] using the SWISS-PROT database[44] (accessed 17/03/2017),

117 and separately using the arthropod and viral subsections of the GenBank nr

118 database (accessed 06/06/2017).

119 ### *Identification of long non-coding RNAs*

120 FEELnc[45] was used for the identification of long non-coding RNAs. The coding

121 transcripts training set was constructed from the 1,047 complete universal single

122 copy orthologous genes found with BUSCO v2.0 (database arthropoda_odb9[42]). The

123 mode "shuffle" was used to generate a training set of lncRNA from the debris of the

124 known coding RNA transcripts.

125 ### *Mapping and differential gene expression analysis*

126 Before mapping, error-corrected raw sequence reads were trimmed using the same

127 parameters as before, but without palindrome trimming used by Trinity. Sequence

128 reads were mapped using Bowtie2[46] V2.2.8, and read counts were calculated using

129 Corset[47] V1.0.6. Differential gene expression was analyzed using DESeq2[48] V1.16.1

130 in RStudio[49] V3.4.1.

131 To reduce the number of sequences for KEGG analysis, the longest contig per

132 cluster was chosen from the combined tissue type and early life-history stage data.

133 The KEGG Automatic Annotation Server (KAAS, http://www.genome.jp/tools/kaas/)

134 was used to generate KEGG pathway maps for each contig using BLAST with the

135 single-directional best hit (SBH) method. All scripts will be deposited on GitHub upon

136 acceptance.

137

138 ### *Statistical analyses*

139 For data analysis, the top 2,000 variably expressed genes across the nine tissue
140 types and the top 500 variably expressed genes across the four larval and four post-
141 larval stages were visualized in a principal component analysis and heatmap using
142 variance-stabilizing transformed read-count data from DESeq2. The gene level
143 dendrograms in the heatmap were created using Pearson's correlation for both the
144 tissue type larval/post-larval stages. Euclidean distance was used to cluster tissue
145 types. All statistical analyses were performed in RStudio. Detailed information on the
146 analyses can be found on GitHub upon acceptance.

147

148 ## **Results**

149 ### *Sequence read data and code availability*

150 In total, nine tissues were sequenced in biological triplicates, as well as pools of
151 eight early life-history stages, resulting in an average of 19.9 M ± 1.6 M (mean ± SD)
152 read pairs per sample and 697 M reads in total (Table 1). After quality trimming,
153 99.5% ± 0.6% (mean ± SD) of reads were retained, indicating a high quality data set
154 (>90% reads with ≥Q30). All read data are available on GenBank through the project
155 ID PRJNA421400.

156 ### *Transcriptome assembly and quality control*

157 The initial combined outputs of all four assemblers comprised of 6,113,055 contigs,
158 which were reduced to 462,772 contigs after filtering with Evidential Gene and
159 combining both "okay" and "alternative" contigs. After clustering with Transfuse, the
160 final assembly consisted of 236,388 transcripts with an assembly size of 226 Mb.
161 These, together with transcript annotations, are available on GenBank. The final
162 transcriptome had a high TransRate score of 0.37, with 88% of all reads successfully
163 mapping back to the transcriptome, and only 3.2% of bases being uncovered. Based
164 on BUSCO, the transcriptome was highly complete with 98% of arthropod ortholog
165 genes being present, and few fragmented or missing genes; however, 51% of the
166 contigs were duplicated/redundant
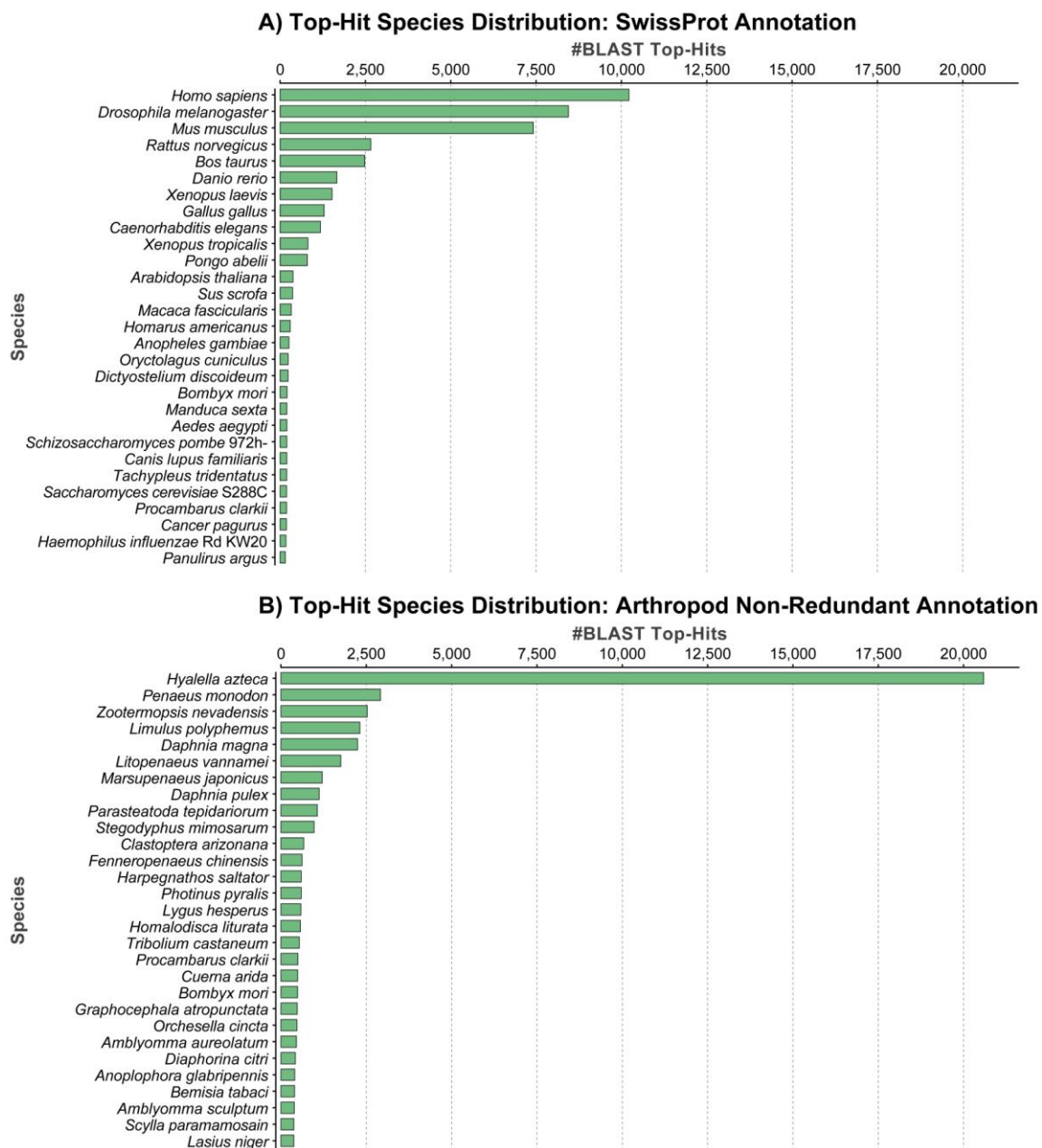167 (C:98.2%[S:46.9%,D:51.3%],F:0.8%,M:1.0%,n:1066).

168 ### *Annotation and gene ontology mapping*

169 Annotation against the SwissProt database using BLASTx resulted in 47,871
170 successfully annotated contigs. Of these, 46,977 were successfully GO mapped, of
171 which 41,069 were completely annotated. The top-hit species distribution was
172 dominated by *Homo sapiens* with over 10,000 hits, followed by *Drosophila*
173 *melanogaster* with just over 8,000 hits; no shrimp species made it into the list
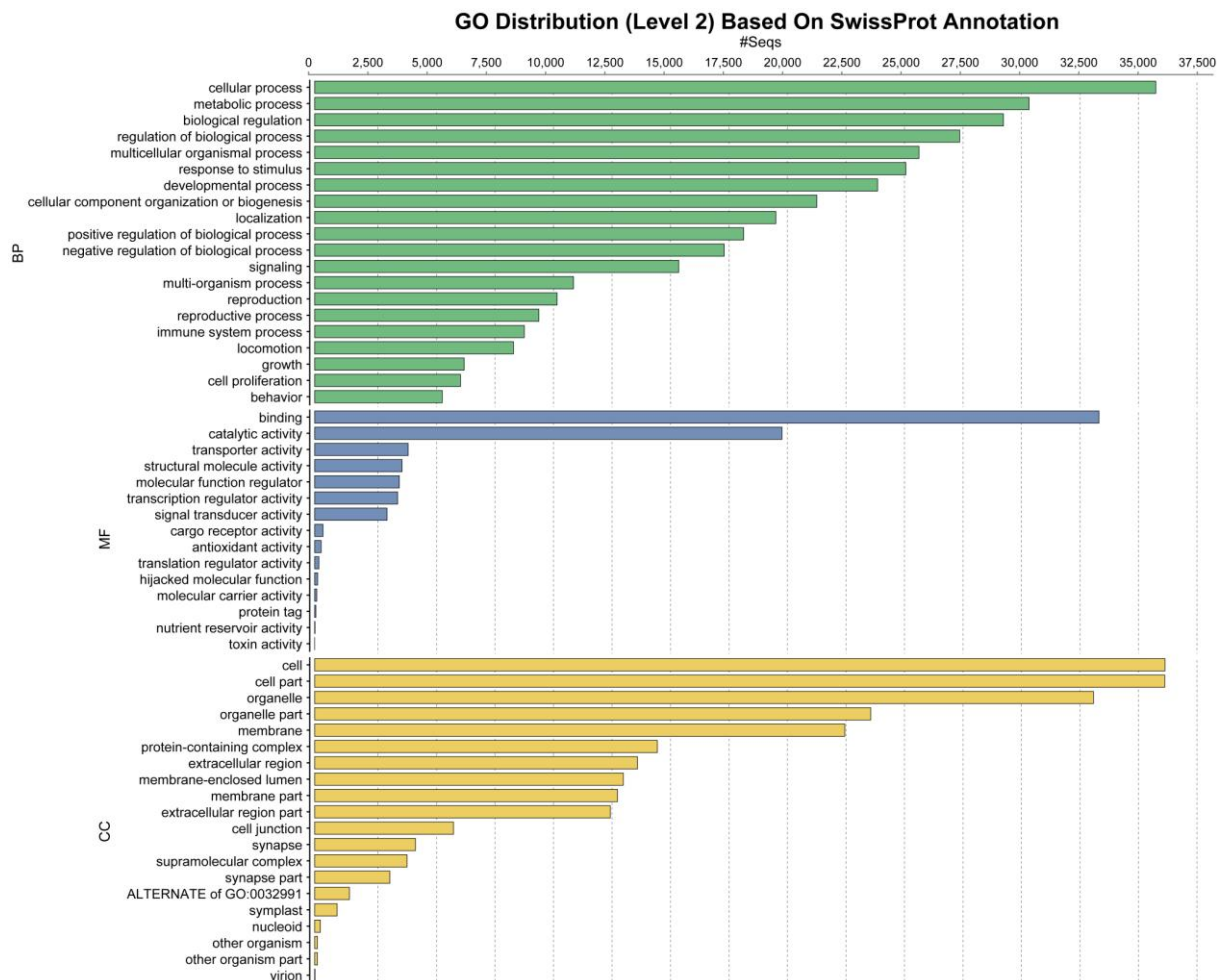
174 (Fig. 1). GO terms for biological processes, molecular function and cellular

175 components were all highly represented in annotated genes (Fig. 2).

176 The annotation against the non-redundant Arthropod (nrA) database using BLASTx

177 resulted in 62,679 successfully annotated contigs, of which 48,456 had a successful

178 GO mapping, and of which 25,201 were completely annotated. The top-hit species

179 distribution was dominated by the freshwater amphipod *Hyalella azteca* with over

180 20,000 hits, followed by *P. monodon* with just over 2,500 hits (Fig. 1). Other penaeid

181 shrimp species included *Litopenaeus vannamei*, *Marsupenaeus japonicus* and

182 *Fenneropenaeus chinensis*, which were the sixth, seventh and twelfth most highly

183 represented species respectively.

184 Detailed information on the annotations can be found in Supplementary Table 1.



185
186 **Figure 1 | Species distribution of successfully annotated sequences across the top 29 species**

187 **using the SwissProt (A) and arthropod subsection of the non-redundant (B) database.**

**GO Distribution (Level 2) Based On SwissProt Annotation**

188

189  **Figure 2 | Distribution of sequence annotations based on the SWISS-PROT database across**

190  **the top 20 GO terms at level 2.** Determined across the entire dataset for Biological Process (BP,

191  green), Molecular Function (MF, blue), and Cellular Component (CC, yellow).

192

193  *Sequence read mapping and differential gene expression analysis*

194  Using Bowtie2, 67.4% ± 4.8% (mean ± SD) of the paired reads successfully mapped

195  to the transcriptome. Using corset for read counting and additional clustering, the

196  initial 236,388 contigs were placed into 99,203 transcript clusters for the nine tissue

197  types and 58,678 transcript clusters for the eight early life-history stages (larval and

198  post-larval stage). A total of 176,966 contigs were used in the clustering of tissues

199  and larvae, with 113,435 shared contigs, 8,188 contigs unique to larvae and 55,343

200  contigs unique to adult tissues.

201  Different tissue types expressed between 9,939 and 12,255 transcript clusters

202  (defined as > 50 normalized read counts per cluster), and between 17 and 316

203  unique sets of transcript clusters (defined as a cluster with > 10 normalized read

204  counts and < 10 normalized read counts in all other tissue types) (Table 3). The

205  ability to annotate transcript clusters varied across tissue types (63.0% to 85.9%). In

206  terms of unique tissue specific transcript clusters, hepatopancreas contained the

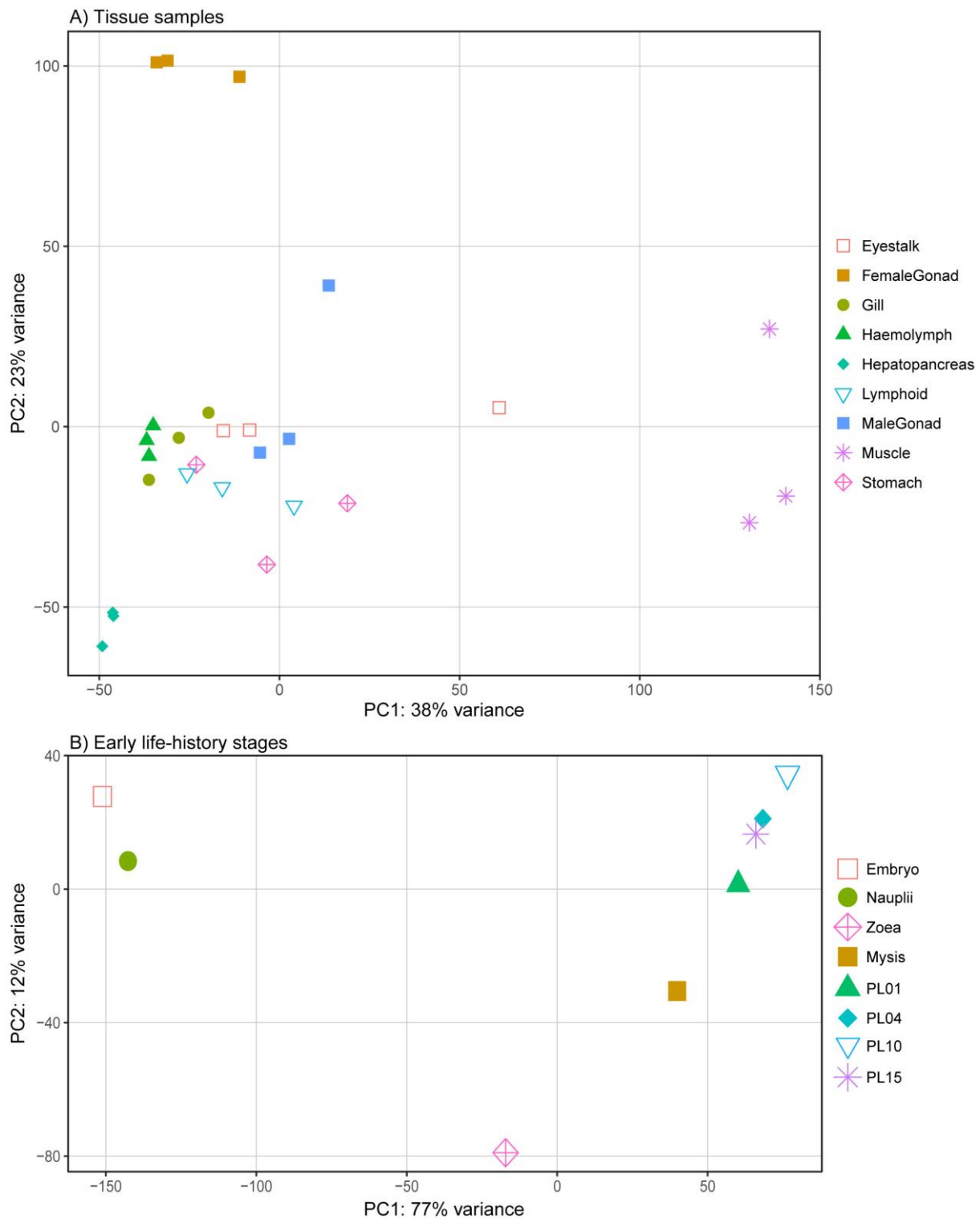207  largest number (316), followed by female gonad (161) and gill (153). Annotation

208    rates of these unique tissue-specific clusters were markedly lower (12.5% to 66.8%)

209    than with clusters shared across all tissue types (82.5% and 85.9%)

210    **Table 3 | Numbers of transcript clusters and cluster annotation rates across**

211    **transcriptomes determined for the nine adult *P. monodon* tissue types analysed.** Total

212    numbers of expressed clusters (>50 normalized read counts), uniquely expressed clusters

213    (normalized read count of >10 in a specific tissue, while having <10 read counts in all other

214    tissues) and constitutively expressed (> 50 normalized read counts in all) clusters within all

215    tissues in this study, and their relative annotation statistics. Numbers represent clusters

216    across all three respective tissue replicates. SP = SWISS-PROT database, nrA = non-

217    redundant Arthropod database.

| Tissue type | Total expressed clusters | | Uniquely expressed clusters | |
|---|---|---|---|---|
| | Number | % Annotated (SP/nrA) | Number | % Annotated (SP/nrA) |
| Eyestalk | 11,173 | 67.3 / 72.8 | 31 | 29.0 / 48.4 |
| Female Gonad | 9,941 | 74.3 / 79.7 | 161 | 37.3 / 45.3 |
| Gill | 12,255 | 63.7 / 69.8 | 153 | 30.7 / 39.2 |
| Haemolymph | 10,577 | 66.1 / 71.4 | 17 | 23.5 / 29.4 |
| Hepatopancreas | 12,169 | 67.7 / 73.9 | 316 | 49.7 / 66.8 |
| Lymphoid Organ | 11,923 | 63.0 / 68.5 | 24 | 54.2 / 66.7 |
| Male Gonad | 10,387 | 71.9 / 77.5 | 71 | 32.4 / 42.3 |
| Muscle | 11,405 | 66.9 / 72.4 | 77 | 33.8 / 48.1 |
| Stomach | 9,939 | 68.6 / 73.7 | 24 | 12.5 / 33.3 |
| Constitutive | 4,300 | 82.5 / 85.9 | - | - |

218

219    A principal component analysis (PCA) of the top 1,000 differentially expressed

220    transcripts across the nine adult tissue types showed strong clustering for most

221    tissue replicates, with the exception of stomach and eyestalk (Fig. 3A).

222    Haemolymph, female gonad and muscle formed distinct clusters separated from

223    other tissues, while eyestalk, gill, haemolymph, lymphoid organ, male gonad and

224    stomach tissues were much more closely associated and showed less distinct

225    clustering (Fig. 3A). A PCA of the top 500 differentially expressed transcripts across

226    the eight early life-history stages showed a strong separation within PC1, with

227    embryo and nauplii segregating substantially from the other early life-history larval

228    stages (Fig. 3B). PC1 explained an extraordinary 77% of the variance in transcript

229    clusters expressed across the different discrete larval stages, which appears to be

230    strongly associated with larval development leading from embryo to post-larval

231    stages.

11

**Figure 3 | Principal component analysis showing the top most highly differentially expressed transcripts of A) nine tissue types (top 1,000) and B) eight early life-history stages (top 500).**
PC = principal component, PL = post-larvae
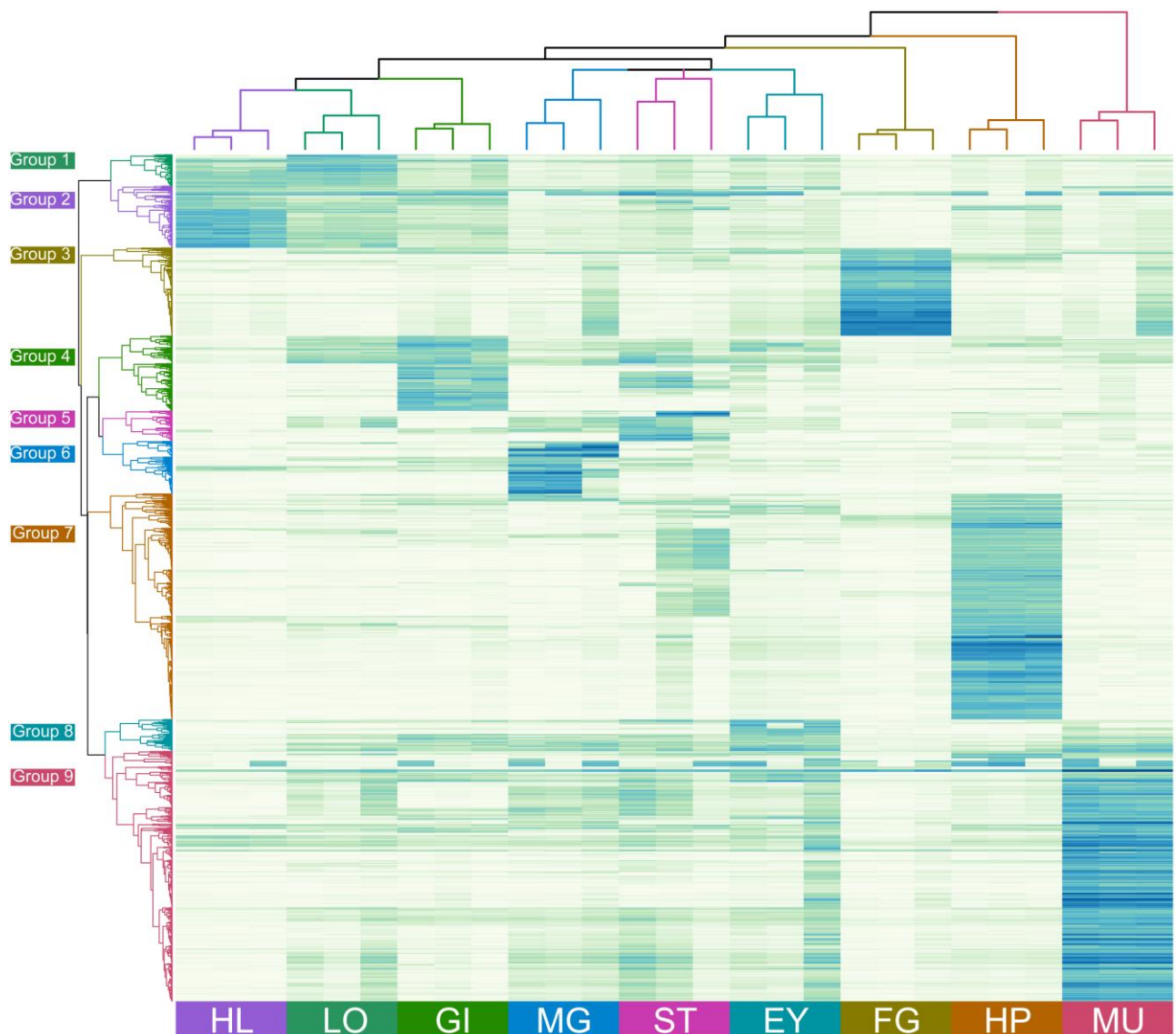
The top 2,000 most variably expressed transcript clusters across all nine tissue types clustered into nine distinct groups using Pearson's correlation (Fig. 4). These groups aligned broadly with expression patterns identified to be unique to each tissues type. For example, group two comprised 208 clusters highly expressed in female gonad, which were mostly successfully annotated (81.8%) using the nrA database.

242   Annotated transcripts included farnesoic acid O-methyltransferase (FAmET),

243   phosphoenolpyruvate carboxykinase (PEPCK), glutathione peroxidase (GPx) and

244   nasrat. Transcripts in each cluster and their annotation are detailed in

245   Supplementary Table 2. Group four consisted of clusters expressed mainly in male

246   gonad that were annotated relatively poorly (38.7%) with many (35.5%) not

247   expressed in the early life-history stages (Table 4). Group nine was the largest and

248   comprised 591 clusters that were mostly annotated (86.0%) and expressed

249   predominantly in muscle tissue. Group seven consisted of 533 clusters that were

250   also mostly annotated (85.7%) and expressed predominantly in hepatopancreatic

251   tissue. Except for male gonad, most clusters expressed in adult tissue types were

252   also expressed in the early life-history stages.

253   **Table 4 | Grouping**s **of the top 2,000 highly variably expressed transcript clusters**
254   **among all nine adult tissue types based on Pearson's correlation.** This includes
255   annotation success and tissue type where each group was predominantly expressed, and the
256   percent of clusters in each group found in adult tissue types but not in the larval stages
257   examined.

| Groups | Predominant tissue type expression site | Number of clusters | % Annotated (SP/nrA) | % in adult but not larval tissues |
|---|---|---|---|---|
| 1 | Lymphoid Organ | 81 | 64.2% / 76.5% | 0.0% |
| 2 | Haemolymph | 139 | 63.3% / 84.9% | 1.4% |
| 3 | Female Gonad | 208 | 55.3% / 81.7% | 6.7% |
| 4 | Gill | 177 | 53.1% / 66.1% | 3.4% |
| 5 | Stomach | 72 | 62.5% / 68.1% | 8.3% |
| 6 | Male Gonad | 124 | 29.0% / 38.7% | 35.5% |
| 7 | Hepatopancreas | 533 | 66.6% / 85.7% | 0.8% |
| 8 | Eyestalk | 75 | 66.7% / 73.3% | 1.3% |
| 9 | Muscle | 591 | 75.1% / 86.0% | 1.5% |
| all | - | 2000 | 64.0% / 84.6% | 4.3% |

258   SP = SWISS-PROT database, nrA = non-redundant Arthropod database

**Figure 4 | Heatmap and hierarchical grouping of the top 2,000 differentially expressed genes in the nine different tissue types.** Gene expression patterns (rows) were grouped into nine expression groups based on Pearson's correlation and the three replicates of each tissue type (columns) into nine tissue groups based on Euclidean distance. EY – eyestalk; FG – female gonad; GI – gill; HL – hemolymph; HP – hepatopancreas; LO – lymphoid organ; MG – male gonad; MU – muscle; ST – stomach.
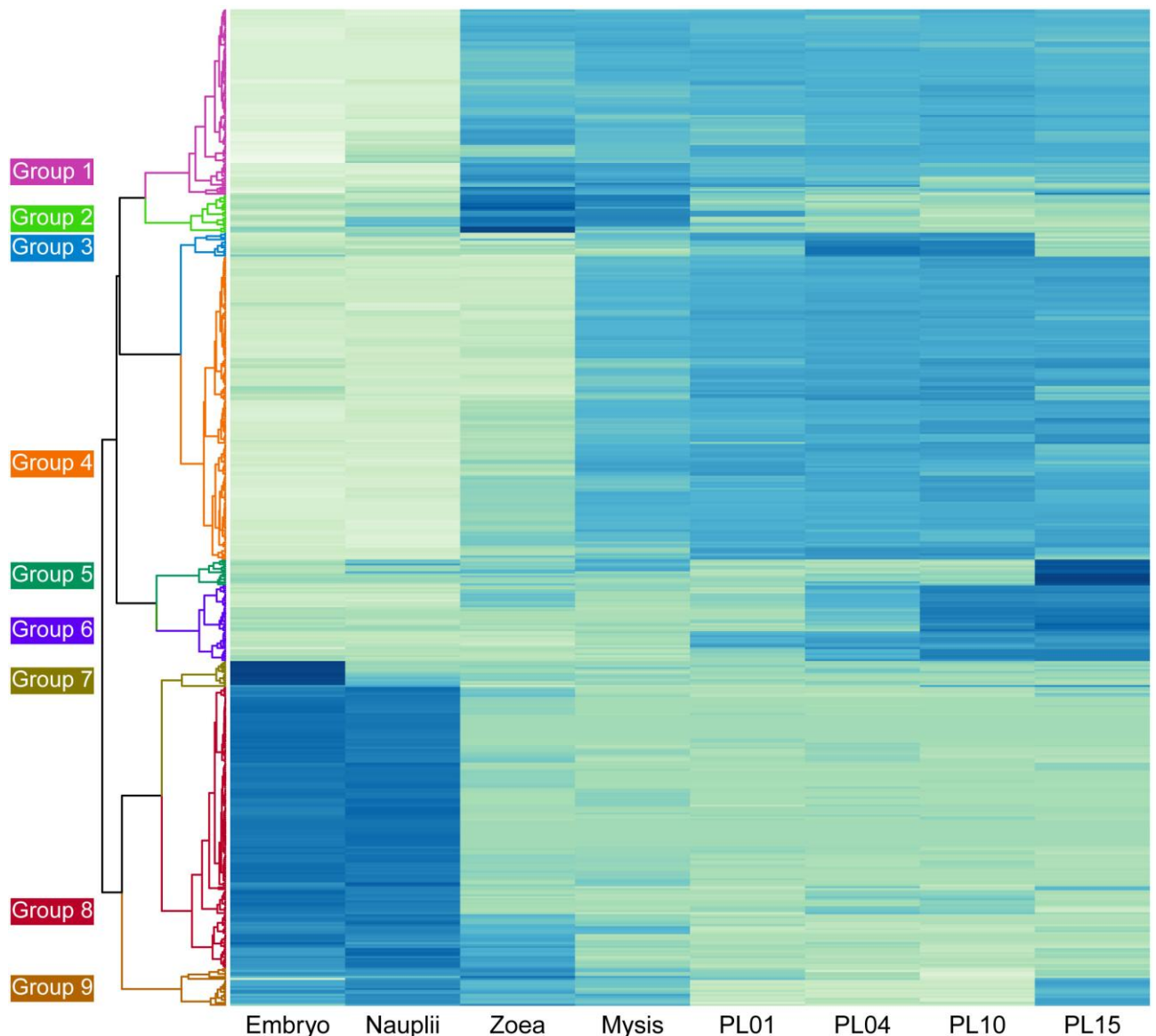
The same top 500 most variably expressed transcript clusters in the different larval and post-larval stages used for the PCA broadly clustered into nine distinct groups based on Pearson's correlation (Fig. 5). Irrespective of the annotation success, the analysis identified transcript clusters that shared similar expression patterns across developmental stages. Embryos and nauplii expressed a set of genes that were not expressed during any other developmental stage (groups 7 and 8). Of the 140 genes expressed exclusively within the embryo and nauplii stages (group 8), only 24.3% and 37.1%, respectively, were annotated successfully using the SWISS-PROT or nrA databases (Table 5). Of the transcript clusters that were annotated, 13 encoded orthologs of the neurotrophic factor *spaetzle* and another 13 encoded orthologs of

14

277 cuticular proteins. Transcripts in each cluster and their annotation are detailed in

278 Supplementary Table 3. Two large clusters of genes were expressed from zoea

279 throughout each subsequent stage (group 1), or from mysis throughout each

280 subsequent stage (group 4). A high percentage (61.2% and 83.1%) of transcripts in

281 these two clusters was annotated. Since each larval stage was sequenced as a pool

282 of individuals, differential gene expression (DGE) analysis could not be performed.

283 **Table 5 | Groupings of the top 500 highly variably expressed transcript clusters among the**
284 **four larval and four post-larval stages based on Pearson's correlation.** This includes annotation
285 success, stages in which transcript groups were predominantly expressed and the percent of clusters
286 in each group found in larval stages, but not in the adult tissue types examined.

| Groups | Stage(s) with predominant expression | Number of clusters | % Annotated (SP/nrA) | % unique to larvae |
|---|---|---|---|---|
| 1 | Mid larval to PL (Z, M, PL01, PL04, PL10, PL15) | 77 | 75.3 / 83.1 | 9.1 |
| 2 | Mid Larval (Z, M) | 35 | 42.9 / 68.6 | 62.9 |
| 3 | Mid PL (PL4, PL10) | 12 | 0.0 / 25.0 | 33.3 |
| 4 | Late larval to PL (M, PL1, PL4, PL10, PL15) | 152 | 61.2 / 69.7 | 18.4 |
| 5 | PL15 | 13 | 69.2 / 92.3 | 76.9 |
| 6 | Late PL (PL4, PL10, PL15) | 38 | 84.2 / 84.2 | 10.5 |
| 7 | Embryo (E) | 12 | 0.0 / 16.7 | 58.3 |
| 8 | Early larval (E, N) | 140 | 24.3 / 37.1 | 85.0 |
| 9 | Larval (E, N, Z, M, PL15) | 21 | 33.3 / 61.9 | 38.1 |
| Total | | 500 | 49.6 / 61.6 | 50.4 |

287 SP = SWISS-PROT database, nrA = non-redundant Arthropod database, E = embryo, N = nauplii, Z =

288 zoea, M = mysis, PL = post larvae (day)

**Fig. 5 | Heatmap and hierarchical grouping of the top 500 differentially expressed genes in the eight larval and post-larval stages examined.** Gene expression patterns in each larval/post-larval stage (row) were grouped into nine expression groups based on Pearson's correlation.

### *Identification of long non-coding RNAs*

We used the set of 1,047 complete USCOs as the training set for classification of coding and non-coding transcripts. It was determined that a coding potential of 0.2642 was the appropriate threshold to balance classification specificity and sensitivity. In total 79,656 transcripts were classified as lncRNAs and the remaining 154,893 transcripts were classified as mRNAs.
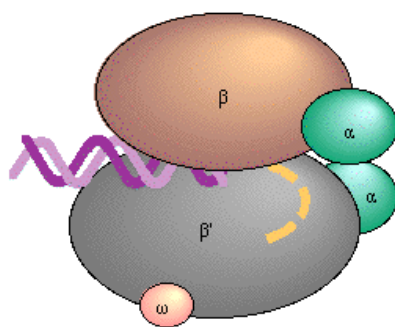
Comparing the lncRNA annotation with the BLASTx annotation, out of the 236,388 contigs 67,960 were uniquely identified as lncRNA, while 13,535 contigs were annotated both as mRNA and lncRNA. At a cluster level, 12,079 out of 58,768 larval clusters (22.6%) and 23,645 out of the 99,203 tissue clusters (23.8%) were uniquely

16

304 annotated as lncRNA. Detailed results of the lncRNA analysis can be found in
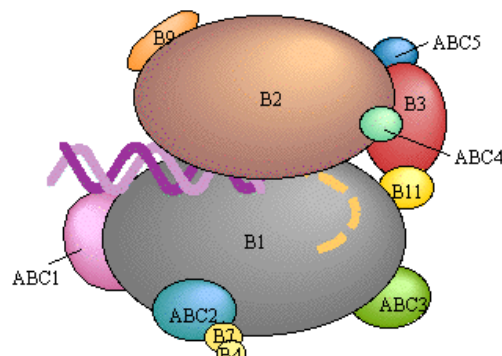
305 Supplementary Table 4.

### *KEGG pathway analysis*

307 Annotated contigs were overlaid onto their respective biological pathways using the

308 Kyoto Encyclopaedia of Genes and Genomes (KEGG) pathways. Genes involved in

309 general eukaryotic cellular processes such as RNA replication (Fig. 6) and basal

310 transcription factor sequences (Fig. 7) were well represented in the *P. monodon*

311 transcriptome. As expected, assignments to KEGG pathways in prokaryotes were

312 rare, as were ribosomal RNA assignments. The various biological processes,

313 metabolism and signalling cascades comprising all 235 KEGG pathways to which

314 transcripts were assigned are detailed in Supplementary Table 5.
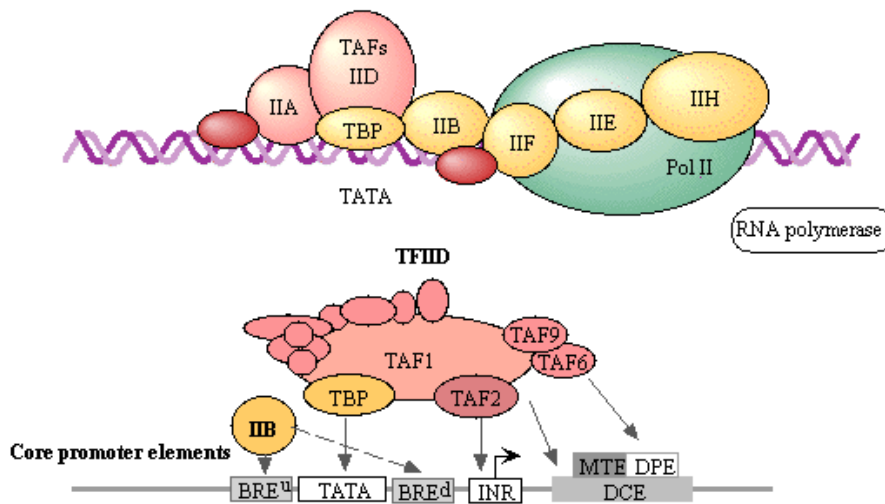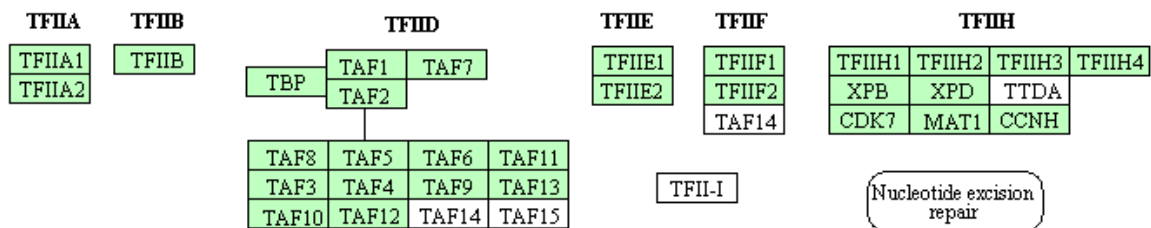


315

**Fig. 6 | Presence of mRNA contigs that encode for RNA polymerase subunits based on KEGG pathway analysis.** Green shading highlights the presence of gene orthologs in the *P. monodon* transcriptome.

319

**Fig. 7 | Presence of eukaryotic basal transcription factor sequences based on KEGG pathway analysis.** Green shading highlights the presence of gene in the *P. monodon* transcriptome.

## *Virus discovery*

Interrogating the *P. monodon* transcriptome against the viral subsection of the non-redundant database using BLASTx assigned viral annotations to 12,744 contigs. Detailed information on the viral blast can be found in Supplementary Table 6. Closer inspection of the data identified the vast majority (>99.8%) of these to represent short motifs conserved between eukaryote cell proteins and related homologs viruses with generally large and complex DNA genomes such as giant viruses, poxviruses, herpes viruses and baculoviruses. Additional BLASTx searches of the GenBank nr database using representative contigs confirmed them to be or likely be endogenous shrimp gene transcripts. The remaining 21 contigs had Top Hit E-value scores identifying them to be related most closely to strains of Gill-associated virus (GAV; 4 contigs, longest 26,235 nt), *Penaeus chinensis* hepandenovirus (*Pchi*HDV; 4 contigs, longest 1,884 nt), Wenzhou shrimp virus 2 (WSV2; RdRp, hypothetical protein and G protein contigs, longest 6,891 nt), Deformed wing virus (DWV; 10,133 nt), Wenzhou shrimp virus 8 (WSV8; 6 contigs,

18

339   longest 4,579 nt), Beihai picorna-like virus 2 (5,277 nt), Wenzhou picorna-like virus

340   23 (551 nt) and Moloney murine leukaemia virus Pr180 sequence (Mo-MuLV; 2,431

341   nt). Lastly, over 1200 contigs with homology to phages were detected, some of

342   which related to phage tail protein and tetracycline resistance.

343

## Discussion

345   Here we report a comprehensive black tiger shrimp (*Penaeus monodon*)

346   transcriptome assembled from nine tissues, four larval stages and four post-larval

347   stages. The transcriptome was generated to expand the genetic resources available

348   for this species to help investigate the genetic basis behind larval developmental

349   stage transitions and tissue functioning, as well as traits with potential to be exploited

350   commercially for the aquaculture of this and other shrimp species. The aim was

351   therefore to generate a highly complete *P. monodon* transcriptome at the risk of it

352   containing higher levels of transcript redundancy. This was confirmed by BUSCO

353   results which demonstrated the transcriptome to be highly complete (C: 98.2%) with

354   low fragmentation (F: 0.8%) or missing (M: 1.0%) genes but high levels of duplication

355   (D: 51.3%). These assembly statistics are comparable to those obtained by a

356   transcriptome assembly from *L. vannamei*[15] (C: 98.0%, F: 0.7%, M: 1.3%, D: 25.5%),

357   but greatly exceeded those of another *P. monodon* assembly focussing on gonadial

358   tissue recently made available publicly[10] (C: 33.7%, F: 44.9%, M: 21.4%, D: 6.8%).

359   As other recent NGS analyses of *P. monodon* have focussed on only one or two

360   tissue types without including any larval stages or biological replicates, generated

361   fewer total reads, or experienced data loss due to quality trimming of low quality

362   reads or low mapping efficiencies[8-11], these are likely to have missed many

363   transcripts. In contrast, the sequencing and assembly strategy used here covered

364   more tissue types at greater read depth and employed multiple *de novo* assembly

365   tools to reduce assembler bias.

### *Functional annotation and comparative analysis*

367   Using the nrA database, 30.0% of transcript clusters found in the nine tissue types

368   and 38.1% of transcript clusters found in the eight larval/post-larval stages analysed

369   were successfully annotated. These annotation levels were comparable to those

370   reported to date in similar studies on different crustaceans[8,15,24,50]. While transcript

371   cluster annotation levels were lower using the SWISS-PROT database compared to

372   the nrA database, the percentage of successful GO-term assignments was

373   substantially higher. In addition to the annotations, analyses were undertaken to

374   identify transcript clusters expressed differentially across tissue types or early life-

375   history stages, irrespective of successful annotation. The identification was done to

376   help provide initial evidence for transcript roles in specific tissue functions or

377   developmental transitions. Despite all efforts made here to improve transcript
378   annotation levels for *P. monodon*, our data reaffirms the need for dedicated
379   functional studies to assign or confirm gene functions of both annotated and
380   unannotated transcript clusters of non-model (crustacean) species.

381   To our best knowledge, to date only two Penaeid shrimp transcriptome assemblies
382   have been made publicly available[10,15], restricting comparative analyses of these
383   transcriptomes. A reciprocal MegaBLAST identified 96.8% of the most recent *P.*
384   *monodon* assembly [10] within the transcriptome described here, but only 40.0% of our
385   assembly was found in the earlier assembly. These comparisons confirm that our
386   transcriptome assembly contains many high quality *P. monodon* transcripts not
387   discovered previously.

388   When compared across species, a reciprocal MegaBLAST showed that the
389   transcriptomes of *P. monodon* (present) and *L. vannamei*[15] shared approximately
390   48% of contigs. Since the assembly metrics of the *L. vannamei* transcriptome were
391   similar to those of our *P. monodon* transcriptome, the low number of shared contigs
392   could stem from considerable differences in transcript type or sequence composition
393   between the two shrimp species. As comprehensive comparisons across crustacean
394   species is currently impractical due to restrictions on publicly-available transcriptome
395   assemblies, the potential value of this warrants effort to consolidate transcriptomic
396   data and to establish both centralized and species-specific databases.

### *Tissue specific expression*

398   Read count data identified independent clusters of transcripts expressed uniquely
399   within different tissues and clusters that formed distinct groups based on their tissue-
400   specific expression patterns. An important consideration for this type of analysis is
401   the normalized read count cutoff value for each cluster to be considered "unique",
402   which was arbitrarily set at above 10 in a specific tissue and < 10 in all others. At
403   >100 normalized read counts, only approximately half of the assigned unique
404   clusters were retained, indicating that the expression levels of many of these
405   potentially tissue-specific clusters was relatively low. Among the annotated transcript
406   clusters most highly expressed in female gonad tissue were FAMeT, PEPCK, GPx
407   and nasrat. Functional roles these proteins may play range from the shrimp moult
408   cycle and reproduction[51], the primary step of gluconeogenesis[52], preventing
409   oxidative stress[52], to specifying terminal regions of the embryo[53]. Among the
410   annotated genes expressed most highly in eyestalk tissue was hyperglycaemic
411   hormone (CHH), a key neuropeptide hormone that regulates blood sugar, moulting
412   and reproduction[54]. A subset of transcript clusters highly expressed in lymphoid
413   organ tissue was also highly expressed in gill tissue, most likely due to high
414   concentrations of haemocytes within both tissue types. The majority of genes
415   expressed most highly in hepatopancreas were annotated, potentially reflecting the

20

416    shared metabolic functions of this organ with those of other animals. Also of much
417    interest were the non-annotated transcripts expressed uniquely in specific tissue
418    types. For example, transcript clusters expressed highly in male gonad were poorly
419    annotated by both databases and included a large proportion of clusters, annotated
420    or not, expressed exclusively in adult tissue types, indicating that male reproductive
421    organs utilize many genes that remain poorly characterized. The grouping of genes
422    with similar expression patterns broadly categorized these transcript clusters into
423    potential functional groups within each tissue type, thereby guiding the selection for
424    more targeted molecular function analyses.

### Larval and post-larval development

426    Based solely on gene expression patterns, the transcriptome data identified unique
427    groups of transcripts involved in transitions between *P. monodon* early life-history
428    stages. There was a major disparity between the annotation success of transcript
429    groups upregulated in early or late stage embryogenesis, highlighting how poorly
430    early developmental pathways have been characterized in crustaceans. Also of
431    significance was the presence of orthologs of the *Spaetzle* gene, known in
432    *Drosophila* flies to establish the dorso-ventral patterning of the early embryo[55] among
433    transcript clusters detected consistently across later larval and post-larval stages.
434    Since each larval and post-larval stage sequenced comprised a pool of several
435    hundred individuals, quantitative and/or spatial transcript expression patterns would
436    be required to draw further functional conclusions. Nevertheless, the data reported
437    here will benefit from similar data on other shrimp and crustacean species,
438    particularly for transcript clusters expressed exclusively in embryo with no significant
439    homology to currently known genes.

### Identification of long non-coding RNAs

441    Long non-coding RNAs (lncRNA) are a type of transcript that have many common
442    features with traditional coding mRNA, including 5' capping, splicing and 3'
443    polyadenylation[56-58]. The nature of lncRNAs is still poorly understood, and it is likely
444    that lncRNAs are in fact a heterogeneous group of transcripts with regulatory
445    functions that are not actively translated into proteins[59]. Thus, their main
446    characteristics are the lack of open reading frames (ORFs) or the presence of non-
447    canonical ORFs in the mature transcript. The biological roles of lncRNAs range from
448    regulation of gene expression, and control of translation, to imprinting. As such, they
449    have been linked to X chromosome inactivation in humans[60], genomic imprinting[61]
450    and cancer[62,63].

451    Due to the lack of a known lncRNA database in shrimp that can be used for their
452    identification, we used FEELnc which scores each transcript according to its coding
453    potential and then selects a threshold score to classify the transcripts into coding or

454  non-coding[45]. This software is particularly useful for non-model species because in
455  the absence of an lncRNA training set, it generates a simulated training set using
456  debris from high confidence coding transcripts. In fly data, this approach showed an
457  MCC value of 0.754 with an accuracy of 0.868[45].

458  In this study, 79,656 transcripts were classified as lncRNAs, of which 67,960 (85.3%)
459  could not be aligned to any protein database. As expected, the use of a non-model
460  organism and the lack of a set with known lncRNA for training led to the ambiguous
461  classification of 13,535 transcripts with low protein-coding potential but clear
462  alignments to known proteins in curated databases. Classification of these
463  transcripts is the first step towards understanding their roles in the development and
464  regulation of gene expression in *Penaeus monodon*.

465  ### *KEGG pathways*

466  Annotated transcript clusters mapped into 235 KEGG pathways (Supplementary
467  Table 3), which have been broadly classified into functional groupings such as
468  general metabolism (e.g. TCA cycle, xenobiotic metabolism, immunity, reproduction),
469  nutritional metabolism (e.g. proteins, lipids, carbohydrates, vitamins), cellular
470  processes (e.g. DNA replication, protein trafficking, apoptosis), biological processes
471  (e.g. circadian rhythm, olfaction and taste, digestion and absorption) and signalling
472  pathways (e.g. PI3K-Akt, MAPK, axis formation, TGF-beta). In general, core
473  pathways such as citrate cycle, oxidative phosphorylation, ribosome biogenesis and
474  RNA/DNA polymerases were better represented than more specific pathways such
475  as the pentose and glucuronate interconversion pathway, or the ascorbate and
476  aldarate metabolism pathway. Furthermore, arthropod specific pathways were
477  generally better represented. For example, the general circadian rhythm pathway
478  was missing several homologs, while the fly specific circadian rhythm pathway was
479  complete. This could be explained by transcripts not sharing sufficient homology with
480  the known genes used for the KEGG analysis and therefore failing to be annotated.
481  Particularly for those pathways highly-conserved among other eukaryotes, the
482  existence of unique transcripts suggests that Penaeid shrimp and possibly
483  crustaceans in general might use metabolic mechanisms differing from eukaryote
484  species studied to date. Their existence also highlights the need for high-quality
485  genome assemblies for shrimp and other crustacean species, overlaid with isoform,
486  tissue-specific and developmental stage transcript expression data, to either help
487  predict gene functions or direct gene knockdown studies, using RNA interference
488  processes as an example, to empirically ascribe functions to novel genes.

489  ### *Virus discovery*

490  Several RNA transcripts and/or genome sequences likely to be from viruses were
491  discovered in the *P. monodon* transcriptome. This was not unexpected considering

492    that it was generated from multiple individuals, tissue types and larval/post-larval

493    stages, as shrimp are co-infected commonly with multiple viruses and as there are

494    several viruses known to be endemic in *P. monodon* populations indigenous to

495    different regions of Australia[64-67]. The presence of near full-length ssRNA genome

496    sequences for viruses such as gill-associated virus (GAV, 26,235 nt) and white spot

497    virus 2 (WSV2, 10,542 nt) provided additional validation of the methods used to

498    synthesize and assemble the transcriptome, and to its completeness as

499    demonstrated by various metrics measuring the nature and number of endogenous

500    gene transcripts. The detection of a ssDNA virus, hepandenovirus, within the

501    transcriptome, presumably detected in a replicative phase, indicates the application

502    of this technique as a tool to also detect the presence of viruses with DNA genomes.

503    In addition to known endemic viruses, the transcriptome contained full-length or near

504    full-length RNA transcripts related closely to the recently-described shrimp viruses

505    WSV2 and WSV8[68,69] unknown until now to occur in Australian *P. monodon*.

506    Moreover, it contained a long transcript (10,133 nt) 95.0% identical to the full- length

507    ssRNA genome of deformed wing virus (DWV), a virus of Varroa mites that is

508    transmitted to honeybees[70], and one of a rapidly expanding number of *Iflavirus*

509    species now being discovered in diverse insect species also including beetles,

510    wasps, caterpillars and moths[71]. As essentially all DWV-like genome sequence reads

511    in this study originated from the stomach of a single individual shrimp, they were

512    potentially derived from a virus-infected honeybee or mite-infested honeybee

513    ingesting by this shrimp. While honeybees infested with Varroa mites have been

514    detected recently in North Queensland not far from where the shrimp was

515    collected[72], DWV itself has not been detected in a comprehensive recent study[73].

516    The present study therefore represents the first detection of a DWV-like genome in

517    Australia, although the origin remains unknown. This reinforces both the strength of

518    the technology in detecting unknown pathogens and also the potential difficulty in

519    interpretation of transcriptome results.

520    A couple of long transcripts of suspected viral origin and expressed across multiple

521    tissue types were also identified. One of these possessed significant BLASTx

522    homology to the reverse transcriptase (RT)-like component of hypothetical protein 1

523    of Beihai picorna-like virus 116 discovered recently in blue swimmer crabs (*Portunus*

524    *pelagicus*)[69]. The other possessed substantial homology to the RT component of the

525    Mo-MuLV Pr180 polyprotein and was expressed across all tissue types except the

526    lymphoid organ, suggesting it to be from a mobile element such as a poly(A)-type

527    retrotransposon or retrovirus[74]. However, determining whether these transcripts

528    containing RT sequences are viral in origin, or represent the products of endogenous

529    retrotransposons like others now being reported in shrimp[75] will require further

530    investigation, as will the nature of the strains, host and distribution ranges,

531 prevalence and potential pathogenicity of the new viruses discovered in the
532 transcriptome.

533 ***Conclusions***

534 This study describes the assembly of a comprehensive and high quality
535 transcriptome from nine different tissue types, and eight larval and post-larval early
536 life-history stages of the black tiger shrimp, *Penaeus monodon*. It also summarizes
537 the number and nature of specific transcript clusters differentially expressed in
538 different tissue types and larval stages, and the Clusters were functionally annotated
539 and mapped to 235 KEGG pathways. Unique transcript clusters and cluster groups
540 were defined across distinct tissues and early life-history stages, providing initial
541 evidence for their roles in specific tissue functions or developmental transitions. The
542 current transcriptome provides a valuable resource for further investigation of
543 directing gene-function studies to increase basic functional biology knowledge in
544 shrimp and for investigating molecular basis of traits of relevance to the aquaculture
545 of shrimp. While the current transcriptome already provides an improved resource for
546 *P. monodon*, further effort is required using long-read sequencing data, such as
547 provided by PacBio, to better resolve genes at isoform level. Lastly, this high-quality
548 *de novo* assembly and data set are publically available and will hopefully support
549 research projects that underpin transformational advances in how we culture shrimp
550 globally.

551

557

558 **References**

559 1    FAO.    Fisheries and Aquaculture topics. The State of World Fisheries and Aquaculture
560      (SOFIA) (Food and Agriculture Organization United Nations, 2016).
561 2    Gjedrem, T., Robinson, N. & Rye, M. The importance of selective breeding in aquaculture to
562      meet future demands for animal protein: a review. *Aquaculture* **350**, 117-129 (2012).
563 3    Jones, D. B. *et al.* A comparative integrated gene-based linkage and locus ordering by
564      linkage disequilibrium map for the Pacific white shrimp, Litopenaeus vannamei. *Scientific*
565      *Reports* **7** (2017).
566 4    Wang, Z., Gerstein, M. & Snyder, M. RNA-Seq: a revolutionary tool for transcriptomics.
567      *Nature reviews genetics* **10**, 57-63 (2009).
568 5    Saha, S. *et al.* Using the transcriptome to annotate the genome. *Nature biotechnology* **20**,
569      508-512 (2002).

570  6   Yu, Y. *et al.* SNP discovery in the transcriptome of white Pacific shrimp Litopenaeus
571      vannamei by next generation sequencing. *PLoS One* **9**, e87218 (2014).

572  7   Song, L., Shankar, D. S. & Florea, L. Rascaf: Improving Genome Assembly with RNA
573      Sequencing Data. *The Plant Genome* **doi:10.3835/plantgenome2016.03.0027** (2016).

574  8   Nguyen, C. *et al.* De novo assembly and transcriptome characterization of major growth-
575      related genes in various tissues of Penaeus monodon. *Aquaculture* **464**, 545-553 (2016).

576  9   Rotllant, G. *et al.* Identification of genes involved in reproduction and lipid pathway
577      metabolism in wild and domesticated shrimps. *Marine genomics* **22**, 55-61 (2015).

578  10  Uengwetwanit, T. *et al.* Transcriptome-based discovery of pathways and genes related to
579      reproduction of the black tiger shrimp (Penaeus monodon). *Marine Genomics* (2017).

580  11  Soonthornchai, W. *et al.* Differentially expressed transcripts in stomach of Penaeus monodon
581      in response to AHPND infection. *Developmental & Comparative Immunology* **65**, 53-63
582      (2016).

583  12  Chen, K. *et al.* Transcriptome and molecular pathway analysis of the hepatopancreas in the
584      Pacific White Shrimp Litopenaeus vannamei under chronic low-salinity stress. *PLoS One* **10**,
585      e0131503 (2015).

586  13  Li, C. *et al.* Analysis of Litopenaeus vannamei transcriptome using the next-generation DNA
587      sequencing technique. *PloS one* **7**, e47442 (2012).

588  14  Chen, X. *et al.* Transcriptome analysis of Litopenaeus vannamei in response to white spot
589      syndrome virus infection. *PLoS One* **8**, e73218 (2013).

590  15  Ghaffari, N. *et al.* Novel transcriptome assembly and improved annotation of the whiteleg
591      shrimp (Litopenaeus vannamei), a dominant crustacean in global seafood mariculture.
592      *Scientific reports* **4**, 7081 (2014).

593  16  Guo, H. *et al.* Trascriptome analysis of the Pacific white shrimp Litopenaeus vannamei
594      exposed to nitrite by RNA-seq. *Fish & shellfish immunology* **35**, 2008-2016 (2013).

595  17  Hu, D., Pan, L., Zhao, Q. & Ren, Q. Transcriptomic response to low salinity stress in gills of
596      the Pacific white shrimp, Litopenaeus vannamei. *Marine genomics* **24**, 297-304 (2015).

597  18  Lu, X. *et al.* Transcriptome analysis of the hepatopancreas in the Pacific white shrimp
598      (Litopenaeus vannamei) under acute ammonia stress. *PloS one* **11**, e0164396 (2016).

599  19  Sookruksawong, S., Sun, F., Liu, Z. & Tassanakajon, A. RNA-Seq analysis reveals genes
600      associated with resistance to Taura syndrome virus (TSV) in the Pacific white shrimp
601      Litopenaeus vannamei. *Developmental & Comparative Immunology* **41**, 523-533 (2013).

602  20  Wei, J., Zhang, X., Yu, Y., Li, F. & Xiang, J. RNA-Seq reveals the dynamic and diverse
603      features of digestive enzymes during early development of Pacific white shrimp Litopenaeus
604      vannamei. *Comparative Biochemistry and Physiology Part D: Genomics and Proteomics* **11**,
605      37-44 (2014).

606  21  Xue, S. *et al.* Sequencing and de novo analysis of the hemocytes transcriptome in
607      Litopenaeus vannamei response to white spot syndrome virus infection. *PLoS One* **8**, e76718
608      (2013).

609  22  Zeng, D. *et al.* Transcriptome analysis of Pacific white shrimp (Litopenaeus vannamei)
610      hepatopancreas in response to Taura syndrome Virus (TSV) experimental infection. *PloS one*
611      **8**, e57515 (2013).

612  23  Zhang, D., Wang, F., Dong, S. & Lu, Y. De novo assembly and transcriptome analysis of
613      osmoregulation in Litopenaeus vannamei under three cultivated conditions with different
614      salinities. *Gene* **578**, 185-193 (2016).

615  24  Powell, D., Knibb, W., Remilton, C. & Elizur, A. De-novo transcriptome analysis of the banana
616      shrimp (Fenneropenaeus merguiensis) and identification of genes associated with
617      reproduction and development. *Marine genomics* **22**, 71-78 (2015).

618  25  Sellars, M. J., Trewin, C., McWilliam, S. M., Glaves, R. & Hertzler, P. L. Transcriptome
619      profiles of Penaeus (Marsupenaeus) japonicus animal and vegetal half-embryos: identification

620          of sex determination, germ line, mesoderm, and other developmental genes. *Marine*
621          *Biotechnology* **17**, 252-265 (2015).

622   26   Powell, D., Knibb, W. & Elizur, A. in *Proceedings of the 24th Plant and Animal Genome*
623          *Conference.* (Plant and Animal Genome (PAG) Conference).

624   27   Powell, D., Knibb, W., Nguyen, N. H. & Elizur, A. Transcriptional profiling of banana shrimp
625          Fenneropenaeus merguiensis with differing levels of viral load. *Integrative and comparative*
626          *biology* **56**, 1131-1143 (2016).

627   28   Wang, W. *et al.* Gill transcriptomes reveal involvement of cytoskeleton remodeling and
628          immune defense in ammonia stress response in the banana shrimp Fenneropenaeus
629          merguiensis. *Fish & shellfish immunology* **71**, 319-328 (2017).

630   29   Li, S., Zhang, X., Sun, Z., Li, F. & Xiang, J. Transcriptome analysis on Chinese shrimp
631          Fenneropenaeus chinensis during WSSV acute infection. *PloS one* **8**, e58627 (2013).

632   30   Shi, X. *et al.* Transcriptome analysis of 'Huanghai No. 2' Fenneropenaeus chinensis response
633          to WSSV using RNA-seq. *Fish & Shellfish Immunology* **75**, 132-138,
634          doi:https://doi.org/10.1016/j.fsi.2018.01.045 (2018).

635   31   Andrews, S. *FastQC: a quality control tool for high throughput sequence data*,
636          <http://www.bioinformatics.babraham.ac.uk/projects/fastqc> (2010).

637   32   MacManes, M. D. Establishing evidenced-based best practice for the de novo assembly and
638          evaluation of transcriptomes from non-model organisms. *bioRxiv*, 035642 (2016).

639   33   Song, L. & Florea, L. Rcorrector: efficient and accurate error correction for Illumina RNA-seq
640          reads. *GigaScience* **4**, 1 (2015).

641   34   Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data without a
642          reference genome. *Nature biotechnology* **29**, 644-652 (2011).

643   35   Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
644          data. *Bioinformatics*, btu170 (2014).

645   36   Liu, J. *et al.* BinPacker: Packing-Based De Novo Transcriptome Assembly from RNA-seq
646          Data. *PLoS Comput Biol* **12**, e1004772 (2016).

647   37   Peng, Y. *et al.* IDBA-tran: a more robust de novo de Bruijn graph assembler for
648          transcriptomes with uneven expression levels. *Bioinformatics* **29**, i326-i334 (2013).

649   38   Chang, Z. *et al.* Bridger: a new framework for de novo transcriptome assembly using RNA-
650          seq data. *Genome biology* **16**, 1 (2015).

651   39   Gilbert, D. *EvidentialGene: tr2aacds, mRNA Transcript Assembly Software*,
652          <http://arthropods.eugenes.org/EvidentialGene/about/EvidentialGene_trassembly_pipe.html>
653          (2013).

654   40   Smith-Unna, R., Boursnell, C., Patro, R., Hibberd, J. & Kelly, S. TransRate: reference free
655          quality assessment of de novo transcriptome assemblies. *Genome research*, gr.
656          196469.196115 (2016).

657   41   Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO:
658          assessing genome assembly and annotation completeness with single-copy orthologs.
659          *Bioinformatics*, btv351 (2015).

660   42   Zdobnov, E. M. *et al.* OrthoDB v9. 1: cataloging evolutionary and functional annotations for
661          animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic acids research* **45**, D744-
662          D749 (2016).

663   43   Conesa, A. *et al.* Blast2GO: a universal tool for annotation, visualization and analysis in
664          functional genomics research. *Bioinformatics* **21**, 3674-3676 (2005).

665   44   Boeckmann, B. *et al.* The SWISS-PROT protein knowledgebase and its supplement TrEMBL
666          in 2003. *Nucleic acids research* **31**, 365-370 (2003).

667   45   Wucher, V. *et al.* FEELnc: a tool for long non-coding RNA annotation and its application to the
668          dog transcriptome. *Nucleic acids research* **45**, e57-e57 (2017).

669  46  Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nature methods*
670      **9**, 357 (2012).

671  47  Davidson, N. M. & Oshlack, A. Corset: enabling differential gene expression analysis for de
672      novo assembled transcriptomes. *Genome biology* **15**, 1 (2014).

673  48  Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
674      RNA-seq data with DESeq2. *Genome biology* **15**, 550 (2014).

675  49  Racine, J. S. RStudio: A Platform-Independent IDE for R and Sweave. *Journal of Applied*
676      *Econometrics* **27**, 167-172 (2012).

677  50  Baranski, M. *et al.* The development of a high density linkage map for black tiger shrimp
678      (Penaeus monodon) based on cSNPs. *PLoS One* **9**, e85413 (2014).

679  51  Homola, E. & Chang, E. S. Methyl farnesoate: crustacean juvenile hormone in search of
680      functions. *Comparative Biochemistry and Physiology Part B: Biochemistry and Molecular*
681      *Biology* **117**, 347-356 (1997).

682  52  Michal, G. & Schomburg, D. *Biochemical pathways: an atlas of biochemistry and molecular*
683      *biology.*  (Wiley New York, 1999).

684  53  Jiménez, G., González-Reyes, A. & Casanova, J. Cell surface proteins Nasrat and Polehole
685      stabilize the Torso-like extracellular determinant in Drosophila oogenesis. *Genes &*
686      *development* **16**, 913-918 (2002).

687  54  Webster, S. G., Keller, R. & Dircksen, H. The CHH-superfamily of multifunctional peptide
688      hormones controlling crustacean metabolism, osmoregulation, moulting, and reproduction.
689      *General and comparative endocrinology* **175**, 217-233 (2012).

690  55  Morisalo, D. & Anderson, K. V. Signaling pathways that establish the dorsal-ventral pattern of
691      the Drosophila embryo. *Annual review of genetics* **29**, 371-399 (1995).

692  56  Chew, G.-L. *et al.* Ribosome profiling reveals resemblance between long non-coding RNAs
693      and 5′ leaders of coding RNAs. *Development* **140**, 2828-2834 (2013).

694  57  Derrien, T. *et al.* The GENCODE v7 catalog of human long noncoding RNAs: analysis of their
695      gene structure, evolution, and expression. *Genome research* **22**, 1775-1789 (2012).

696  58  Guttman, M., Russell, P., Ingolia, N. T., Weissman, J. S. & Lander, E. S. Ribosome profiling
697      provides evidence that large noncoding RNAs do not encode proteins. *Cell* **154**, 240-251
698      (2013).

699  59  Engreitz, J. M., Ollikainen, N. & Guttman, M. Long non-coding RNAs: spatial amplifiers that
700      control nuclear structure and gene expression. *Nature Reviews Molecular Cell Biology* **17**,
701      756 (2016).

702  60  Brockdorff, N. *et al.* The product of the mouse Xist gene is a 15 kb inactive X-specific
703      transcript containing no conserved ORF and located in the nucleus. *Cell* **71**, 515-526 (1992).

704  61  Koerner, M. V., Pauler, F. M., Huang, R. & Barlow, D. P. The function of non-coding RNAs in
705      genomic imprinting. *Development* **136**, 1771-1783 (2009).

706  62  Leucci, E. *et al.* Melanoma addiction to the long non-coding RNA SAMMSON. *Nature* **531**,
707      518 (2016).

708  63  Prensner, J. R. & Chinnaiyan, A. M. The emergence of lncRNAs in cancer biology. *Cancer*
709      *discovery* **1**, 391-407 (2011).

710  64  Cowley, J. A., Dimmock, C. M., Spann, K. M. & Walker, P. J. Gill-associated virus of Penaeus
711      monodon prawns: an invertebrate virus with ORF1a and ORF1b genes related to arteri-and
712      coronaviruses. *Journal of General Virology* **81**, 1473-1484 (2000).

713  65  Cowley, J. A. *et al. Tactical Research Fund: Aquatic Animal Health Subprogram: Viral*
714      *presence, prevalence and disease management in wild populations of the Australian Black*
715      *Tiger prawn (Penaeus monodon).*  (FRDC, 2015).

716  66  Mohr, P. G. *et al.* New yellow head virus genotype (YHV7) in giant tiger shrimp Penaeus
717      monodon indigenous to northern Australia. *Diseases of aquatic organisms* **115**, 263-268
718      (2015).

719  67  Owens, L., La Fauce, K. & Claydon, K. The effect of Penaeus merguiensis densovirus on
720      Penaeus merguiensis production in Queensland, Australia. *Journal of fish diseases* **34**, 509-
721      515 (2011).

722  68  Li, C.-X. *et al.* Unprecedented genomic diversity of RNA viruses in arthropods reveals the
723      ancestry of negative-sense RNA viruses. *Elife* **4** (2015).

724  69  Shi, M. *et al.* Redefining the invertebrate RNA virosphere. *Nature* **540**, 539 (2016).

725  70  Lanzi, G. *et al.* Molecular and biological characterization of deformed wing virus of honeybees
726      (Apis mellifera L.). *Journal of virology* **80**, 4998-5009 (2006).

727  71  Valles, S. *et al.* ICTV virus taxonomy profile: Iflaviridae. *Journal of General Virology* **98**, 527-
728      528 (2017).

729  72  Department of Agriculture and Fisheries, Q. G. *Varroa mite detection in Townsville*,
730      <https://www.daf.qld.gov.au/business-priorities/animal-industries/bees/diseases-and-
731      pests/asian-honey-bees/general-information-on-varroa-mites> (2017).

732  73  Roberts, J. M., Anderson, D. L. & Durr, P. A. Absence of deformed wing virus and Varroa
733      destructor in Australia provides unique perspectives on honeybee viral landscapes and colony
734      losses. *Scientific Reports* **7**, 6925 (2017).

735  74  Shen, C.-H. & Steiner, L. A. Genome structure and thymic expression of an endogenous
736      retrovirus in zebrafish. *Journal of virology* **78**, 899-911 (2004).

737  75  Sakaew, W., Pratoomthai, B., Pongtippatee, P., Flegel, T. W. & Withyachumnarnkul, B.
738      Discovery and partial characterization of a non-LTR retrotransposon that may be associated
739      with abdominal segment deformity disease (ASDD) in the whiteleg shrimp Penaeus
740      (Litopenaeus) vannamei. *BMC veterinary research* **9**, 189 (2013).

741

## Supplementary Material

743  **Supplementary Table 1** contains all annotation results from the blast against
744  SwissProt and nrA (arthropod subsection of the nr database), including blast metrics,
745  GO terms, interpro scan results and simplified lncRNA results.

746  **Supplementary Table 2** contains the top 2,000 differentially expressed genes in the
747  nine different tissue types used in the heatmap. The table shows normalised
748  expression values for each sample, the nine groupings based on Pearson's
749  correlation presented in the heatmap, and the associated SwissProt and nrA
750  annotations.

751  **Supplementary Table 3** contains the top 500 differentially expressed genes in the
752  eight early life-history stages used in the heatmap. The table shows normalised
753  expression values for each sample, the nine groupings based on Pearson's
754  correlation presented in the heatmap, and the associated SwissProt and nrA
755  annotations.

756  **Supplementary Table 4** contains the detailed results of the lncRNA analysis using
757  the FEELnc pipeline.

758  **Supplementary Table 5** contains links to the 235 KEGG pathway figures based on
759  the transcriptome generated in this study.

760  **Supplementary Table 6** contains all successful blast hits against the viral
761  subsection of the nr database.